

Scalable Semidefinite Programming*

Alp Yurtsever[†], Joel A. Tropp[‡], Olivier Fercoq[§], Madeleine Udell[¶], and Volkan Cevher[†]

Abstract. Semidefinite programming (SDP) is a powerful framework from convex optimization that has striking potential for data science applications. This paper develops a provably correct algorithm for solving large SDP problems by economizing on both the storage and the arithmetic costs. Numerical evidence shows that the method is effective for a range of applications, including relaxations of **MaxCut**, abstract phase retrieval, and quadratic assignment. Running on a laptop, the algorithm can handle SDP instances where the matrix variable has over 10^{13} entries.

Key words. Augmented Lagrangian, conditional gradient method, convex optimization, dimension reduction, first-order method, Frank–Wolfe method, MaxCut, phase retrieval, quadratic assignment, randomized linear algebra, semidefinite programming, sketching.

AMS subject classifications. Primary: 90C22, 65K05. Secondary: 65F99.

1. Motivation. For a spectrum of challenges in data science, methodologies based on semidefinite programming offer remarkable performance both in theory and for small problem instances. Even so, practitioners often critique this approach by asserting that it is impossible to solve semidefinite programs (SDPs) at the scale demanded by real-world applications. We would like to argue against this article of conventional wisdom.

This paper proposes a new algorithm, called **SketchyCGAL**, that can solve very large SDPs to moderate accuracy. The algorithm marries a primal–dual optimization technique to a randomized sketch for low-rank matrix approximation. For every standard-form SDP that satisfies strong duality, **SketchyCGAL** provably finds a near-optimal low-rank approximation of a solution using limited storage and arithmetic. In contrast, given the same computational resources, other methods for large-scale SDP may fail where **SketchyCGAL** succeeds. In particular, **SketchyCGAL** needs dramatically less storage than Burer–Monteiro factorization [27, 22] for difficult problem instances [105].

In addition to the strong theoretical guarantees, we offer extensive evidence that **SketchyCGAL** is a practical optimization algorithm. For example, on a laptop, we can solve the

*Submitted to the editors 6 December 2019.

Funding: VC and AY have received funding from the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation programme under the grant agreement number 725594 (time-data) and the Swiss National Science Foundation (SNSF) under the grant number 200021_178865/1. JAT gratefully acknowledges ONR Awards N00014-11-1-0025, N00014-17-1-2146, and N00014-18-1-2363. MU gratefully acknowledges DARPA Award FA8750-17-2-0101.

[†]Laboratory for Information and Inference Systems, Department of Electrical Engineering, École Polytechnique Fédérale de Lausanne, Lausanne, Switzerland (alp.yurtsever@epfl.ch, volkan.cevher@epfl.ch, <https://lions.epfl.ch/>).

[‡]Department of Computing and Mathematical Sciences, California Institute of Technology, Pasadena, CA, USA (jtropp@cms.caltech.edu, <http://users.cms.caltech.edu/~jtropp>).

[§]Laboratoire Traitement et Communication d'Information, Télécom Paris, Institut Polytechnique de Paris, Palaiseau, France (olivier.fercoq@telecom-paris.fr, <https://perso.telecom-paristech.fr/ofercoq/>).

[¶]Department of Operations Research and Information Engineering, Cornell University, Ithaca, NY, USA (udell@cornell.edu, <https://people.orie.cornell.edu/mru8/>).

MaxCut SDP for a sparse graph with almost 8 million vertices, where the matrix variable has about 10^{13} entries. We also tackle large phase retrieval problems arising from Fourier ptychography [55], as well as relaxations [117, 26] of the quadratic assignment problem.

We will first explain how our approach applies to the MaxCut SDP, and then we will enlarge our scope to include all standard-form SDPs.

1.1. Example: The maximum cut in a graph. To begin, we derive a fundamental SDP [37, 38, 47] that arises in combinatorial optimization. This example shows how SDPs arise from matrix optimization problems with rank constraints. It highlights why large SDPs are computationally challenging, and it allows us to illustrate the potential of our approach.

1.1.1. MaxCut. Consider an undirected graph $G = (V, E)$ composed of the vertex set $V = \{1, \dots, n\}$ and a set E that lists the m edges. The combinatorial Laplacian of the graph is the real positive-semidefinite (psd) matrix

$$(1.1) \quad L := \sum_{\{i,j\} \in E} (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^* \in \mathbb{R}^{n \times n},$$

where $\mathbf{e}_i \in \mathbb{R}^n$ denotes the i th standard basis vector and $*$ refers to the (conjugate) transpose of a matrix or vector. The Laplacian L has at most $2m + n$ nonzero entries.

A *cut* is a subset S of the vertex set V , and the signed indicator of the cut is the vector

$$\chi_S \in \mathbb{R}^n \quad \text{where} \quad \chi_S(i) = \begin{cases} +1, & i \in S; \\ -1, & i \in V \setminus S. \end{cases}$$

Define $\text{weight}(S)$ to be the number of edges with one vertex in S and the other in $V \setminus S$. We can compute the weight of a cut algebraically: $4 \text{weight}(S) = \chi_S^* L \chi_S$. As a consequence, we can search for the maximum-weight cut in the graph via the discrete optimization problem

$$(1.2) \quad \text{maximize} \quad \chi^* L \chi \quad \text{subject to} \quad \chi \in \{\pm 1\}^n.$$

The formulation (1.2) is NP-hard [63]. One remedy is to relax it to an SDP.

Consider the matrix $\mathbf{X} = \chi \chi^*$ where $\chi \in \{\pm 1\}^n$. The matrix \mathbf{X} is psd; its diagonal entries equal one; and it has rank one. We can express the MaxCut problem (1.2) in terms of the matrix \mathbf{X} by rewriting the objective as a trace. Bringing forward the implicit constraints on \mathbf{X} and dropping the rank constraint, we arrive at the MaxCut SDP:

$$(1.3) \quad \text{maximize} \quad \text{tr}(L\mathbf{X}) \quad \text{subject to} \quad \text{diag}(\mathbf{X}) = \mathbf{1}, \quad \mathbf{X} \text{ is psd.}$$

As usual, diag extracts the diagonal of a matrix, and $\mathbf{1} \in \mathbb{R}^n$ is the vector of ones.

The matrix solution \mathbf{X}_* of (1.3) does not immediately yield a cut. Let $\mathbf{x}_* \mathbf{x}_*^*$ be a best rank-one approximation of \mathbf{X}_* with respect to the Frobenius norm. Then the vector $\chi_* = \text{sgn}(\mathbf{x}_*)$ is the signed indicator of a cut S_* . In many cases, the cut S_* yields an excellent solution to the MaxCut problem (1.2). We can also use \mathbf{X}_* to compute a cut that is provably near-optimal via a more involved randomized rounding procedure [47].

1.1.2. What's the issue? We specify an instance of the MaxCut SDP (1.3) by means of the Laplacian \mathbf{L} of the graph, which has $O(m+n)$ nonzero entries. Our goal is to compute a best rank-one approximation of a solution to the SDP, which has $O(n)$ degrees of freedom. In other words, the total cost of representing the input and output of the problem is $O(m+n)$. Sadly, the matrix variable in (1.3) seems to require storage $O(n^2)$. For example, a graph G with two million vertices leads to an SDP (1.3) with almost two trillion real variables!

Storage is one of the main reasons that it has been challenging to solve large instances of the MaxCut SDP reliably. Undeterred, we raise a question:

Can we provably find a best rank-one approximation of a solution to the MaxCut SDP (1.3) with working storage $O(m+n)$? Can we achieve storage $O(n)$?

We are not aware of any correct algorithm that can solve an arbitrary instance of (1.3) with a working storage guarantee better than $\Theta(\min\{m, n^{3/2}\})$; see section 8.

In addition to the limit on storage, a good algorithm should interact with the Laplacian \mathbf{L} only through noninvasive, low-cost operations, such as matrix–vector multiplication.

1.1.3. A storage-optimal algorithm for the MaxCut SDP. Surprisingly, it is possible to achieve all the goals announced in the last paragraph.

Informal Theorem 1.1 (MaxCut via SketchyCGAL). *For any $\varepsilon, \zeta > 0$ and any Laplacian \mathbf{L} , the SketchyCGAL algorithm computes a $(1+\zeta)$ -optimal rank-one approximation of an ε -optimal point of (1.3); see subsection 2.2. The working storage is $O(n/\zeta)$. The algorithm performs at most $\tilde{O}(\varepsilon^{-3})$ matrix–vector multiplies with the Laplacian \mathbf{L} , plus lower-order arithmetic. The algorithm is randomized; it succeeds with high probability over its random choices.*

Informal Theorem 1.1 follows from Theorem 6.3. As usual, \tilde{O} suppresses constants and logarithmic factors. In our experience, SketchyCGAL works *better* than the theorem suggests: we find that SketchyCGAL succeeds using only $\tilde{O}(\varepsilon^{-5/4})$ matrix–vector multiplies with \mathbf{L} , provided the slightly larger storage budget $\tilde{O}((\varepsilon^{-1/4} + \zeta^{-1})n)$. See subsection 7.2 for numerical examples on sparse graphs with almost 10^7 vertices.

Remark 1.2 (Prior work). In contrast to Burer–Monteiro methods [105] and to the approximate complementarity paradigm [39], SketchyCGAL succeeds for every instance of MaxCut. At present, the fastest theoretical algorithm [69] for finding an ε -optimal solution to (1.3) requires $O(m)$ working storage and $\tilde{O}(\varepsilon^{-3.5})$ matrix–vector multiplies with \mathbf{L} .

1.2. A model problem. SketchyCGAL can solve all standard-form SDPs. To simplify some aspects of the presentation, we will focus on a model problem that includes an extra trace constraint. Appendix D explains how to extend SketchyCGAL to a more expressive problem template that includes standard-form SDPs with additional (conic) inequality constraints.

1.2.1. The trace-constrained SDP. We work over the field $\mathbb{F} = \mathbb{R}$ or $\mathbb{F} = \mathbb{C}$. For each $n \in \mathbb{N}$, define the set $\mathbb{S}_n := \mathbb{S}_n(\mathbb{F})$ of (conjugate) symmetric $n \times n$ matrices with entries in \mathbb{F} .

Introduce the set of psd matrices with trace one:

$$(1.4) \quad \Delta_n := \{\mathbf{X} \in \mathbb{S}_n : \text{tr } \mathbf{X} = 1 \text{ and } \mathbf{X} \text{ is psd}\}.$$

Our model problem is the following trace-constrained SDP:

$$(1.5) \quad \begin{aligned} & \text{minimize} && \text{tr}(\mathbf{C}\mathbf{X}) \\ & \text{subject to} && \text{tr}(\mathbf{A}_i\mathbf{X}) = b_i \quad \text{for } i = 1, \dots, d \quad \text{and} \quad \mathbf{X} \in \alpha\Delta_n. \end{aligned}$$

The trace parameter $\alpha > 0$, each matrix $\mathbf{C}, \mathbf{A}_1, \dots, \mathbf{A}_d \in \mathbb{S}_n$, and $b_1, \dots, b_d \in \mathbb{R}$. We always assume that (1.5) satisfies strong duality with its standard-form dual problem.

Remark 1.3 (Standard-form SDPs). To solve a standard-form SDP, we replace the inclusion $\mathbf{X} \in \alpha\Delta_n$ with the constraints that \mathbf{X} is psd and $\text{tr } \mathbf{X} \leq \alpha$ for a large enough parameter α .

1.2.2. Applications. The model problem (1.5) has diverse applications in statistics, signal processing, quantum information theory, combinatorics, and beyond. Evidently, the MaxCut SDP (1.3) is a special case of the model problem (1.5). The template (1.5) also covers alignment problems that appear in biological imaging [17] and in robotics [92]. It includes SDPs that arise from phase retrieval problems in imaging sciences [33, 104, 55] and in segmentation problems from computer vision [56]. It also supports contemporary machine learning tasks, such as certifying robustness of neural networks [90]. The possibilities are endless.

1.3. Complexity of SDP formulations and solutions. This section describes some special features that commonly appear in large, real-world SDPs. Our algorithm will exploit these features in an essential way, even as it provides guarantees for every instance of (1.5).

1.3.1. Structure of the problem data. The matrices \mathbf{C} and \mathbf{A}_i that appear in (1.5) are often highly structured, or sparse, or have low-rank. As such, we can specify the SDP using a small amount of information. In our work, we exploit this property by treating the problem data for the SDP (1.5) as a collection of black boxes that support specific linear algebraic operations. The algorithm for solving the SDP only needs to access the data via these black boxes, and we can make sure that these subroutines are implemented efficiently.

1.3.2. Low-rank solutions of SDPs. We will also capitalize on the fact that SDPs frequently have low-rank solutions, or the solutions are approximated well by low-rank matrices. There are several reasons why we can make this surmise.

Weakly-constrained SDPs. First, many SDPs have low-rank solutions just because they are weakly constrained. That is, the number d of linear inequalities in (1.5) is much smaller than the $\Theta(n^2)$ number of entries in the matrix variable. This situation often occurs in signal processing and statistics problems, where d reflects the amount of measured data.

Weakly constrained SDPs provably have low-rank solutions because of the geometry of the set of psd matrices; see [18, Prop. II(13.4) and Prob. II.14.5] and [87].

Fact 1.4 (Barvinok–Pataki). *When $\mathbb{F} = \mathbb{R}$, the SDP (1.5) has a solution with rank $r \leq \sqrt{2(d+1)}$. When $\mathbb{F} = \mathbb{C}$, there is a solution with rank $r \leq \sqrt{d+1}$.*

For example, the MaxCut SDP (1.3) admits a solution with rank $\sqrt{2(n+1)}$.

Although a weakly-constrained SDP can have solutions with high rank, a *generic* weakly-constrained SDP only admits a low-rank solution [5].

Fact 1.5 (Alizadeh et al.). *Let $\mathbb{F} = \mathbb{R}$. Except for a set of matrices $\{\mathbf{C}, \mathbf{A}_1, \dots, \mathbf{A}_d\}$ with measure zero, the solution set of the SDP (1.5) is a unique matrix with rank $r \leq \sqrt{2(d+1)}$.*

Our algorithm targets weakly constrained SDPs, but it works for all instances of (1.5).

Matrix rank minimization. Second, some SDPs are *designed* to produce a low-rank matrix that satisfies a system of linear matrix equations. This idea goes back to the control literature [76, 86], and it was explored thoroughly in Fazel’s thesis [40]. Early applications include Euclidean distance matrix completion [2] and collaborative filtering [96]. Extensive empirical work indicates that these SDPs often produce low-rank solutions.

Structural properties. There are other reasons that an SDP must have a low-rank solution. For instance, consider the optimal power flow SDPs developed by Lavaei and Low [68], where the rank of the solution is controlled by the geometry of the power grid.

1.3.3. Algorithms? To summarize, many realistic SDPs have structured data, and they admit solutions that are (close to) low rank. Are there algorithms that can exploit these features? Although there are a number of methods that attempt to do so, none can provably solve every SDP while controlling storage and arithmetic costs. See [section 8](#) for related work.

1.3.4. Finding low-rank solutions? Why has it been so difficult to develop provably correct algorithms for finding low-rank solutions to structured SDPs? Most approaches that try to control the rank run headlong into a computational complexity barrier: For any fixed rank parameter r , it is NP-hard to solve the model problem (1.5) if the variable \mathbf{X} is also constrained to be a rank- r matrix [40, p. 7].

To escape this sticky fact, we revise the computational goal, following [115]. The key insight is to seek a rank- r matrix that *approximates a solution* to (1.5). See [subsection 2.2](#) for a detailed explanation. This shift in perspective opens up new algorithmic prospects.

1.4. Contributions. This paper explains how to harness the favorable properties that are common in large SDPs. The SketchyCGAL algorithm solves the SDP using a primal–dual optimization method [110] developed by a subset of the authors. Instead of storing the psd matrix variable, we maintain a compressed representation by means of a matrix sketching technique [100]. After the optimization algorithm terminates, we extract from the sketch a low-rank approximation of the solution of the SDP. This idea leads to a practical, provably correct SDP solver that economizes on storage and arithmetic.

Informal Theorem 1.6 (The model problem via SketchyCGAL). *Assume that the model problem (1.5) satisfies strong duality. For any $\varepsilon, \zeta > 0$ and any rank parameter r , the SketchyCGAL algorithm computes a $(1 + \zeta)$ -optimal rank- r approximation of an ε -optimal point of (1.5); see [subsection 2.2](#). The storage cost is $O(d + rn/\zeta)$. Most of the arithmetic consists of $\tilde{O}(\varepsilon^{-3})$ matrix–vector products with each matrix $\mathbf{C}, \mathbf{A}_1, \dots, \mathbf{A}_d$ from the problem data. The algorithm succeeds with high probability.*

[Theorem 6.3](#) contains full theoretical details. Note that the storage $O(d + rn)$ is the minimum for any primal–dual algorithm that returns a rank- r solution to (1.5). The arithmetic requirements can be reduced at the expense of additional storage. Computational evidence indicates that the true arithmetic costs are much lower than we can prove; see [section 7](#). We believe SketchyCGAL is the first SDP solver with these guarantees. And it actually works!

1.5. Roadmap. [Section 2](#) presents an abstract framework for studying SDPs that highlights the challenges associated with large problems. [Section 3](#) outlines a primal–dual algorithm, called CGAL, for solving the model problem (1.5). [Sections 4](#) and [5](#) introduce some

methods from randomized linear algebra that we use to control storage and arithmetic costs. [Section 6](#) develops the SketchyCGAL algorithm, its convergence theory, and its resource usage guarantees. [Appendix D](#) describes some extensions of the model problem. We present a numerical study of SketchyCGAL in [section 7](#). [Section 8](#) discusses related work.

1.6. Notation. The symbol $\|\cdot\|_F$ denotes the Frobenius norm, while $\|\cdot\|_*$ is the nuclear norm (i.e., Schatten-1). The unadorned norm $\|\cdot\|$ refers to the ℓ_2 norm of a vector, the spectral norm of a matrix, or the operator norm of a linear map from $(\mathbb{S}_n, \|\cdot\|_F)$ to $(\mathbb{R}^d, \|\cdot\|)$. We write $\langle \cdot, \cdot \rangle$ for both the ℓ_2 inner product on vectors and the trace inner product on matrices.

The map $\llbracket \mathbf{M} \rrbracket_r$ returns an r -truncated singular-value decomposition of the matrix \mathbf{M} , which is a best rank- r approximation with respect to every unitarily invariant norm [77].

We use the standard computer science interpretation of the asymptotic notation O, \tilde{O}, Θ .

2. Scalable semidefinite programming. To solve the model problem (1.5) efficiently, we need to exploit structure inherent in the problem data. This section outlines an abstract approach that directs our attention to the core computational difficulties.

2.1. Abstract form of the model problem. Let us instate compact notation for the linear constraints in the model problem (1.5). Define a linear map \mathcal{A} and its adjoint \mathcal{A}^* via

$$(2.1) \quad \begin{aligned} \mathcal{A} : \mathbb{S}_n &\rightarrow \mathbb{R}^d \quad \text{where} \quad \mathcal{A}\mathbf{X} = [\langle \mathbf{A}_1, \mathbf{X} \rangle \quad \dots \quad \langle \mathbf{A}_d, \mathbf{X} \rangle]; \\ \mathcal{A}^* : \mathbb{R}^d &\rightarrow \mathbb{S}_n \quad \text{where} \quad \mathcal{A}^*\mathbf{z} = \sum_{i=1}^d z_i \mathbf{A}_i. \end{aligned}$$

We equip the linear map with the operator norm $\|\mathcal{A}\| := \|\mathcal{A}\|_{F \rightarrow \ell_2}$. Form the vector $\mathbf{b} := (b_1, \dots, b_d) \in \mathbb{R}^d$ of constraint values. In this notation, (1.5) becomes

$$(2.2) \quad \text{minimize} \quad \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{subject to} \quad \mathcal{A}\mathbf{X} = \mathbf{b}, \quad \mathbf{X} \in \alpha \Delta_n.$$

Problem instances are parameterized by the tuple $(\mathbf{C}, \mathcal{A}, \mathbf{b}, \alpha)$.

2.2. Approximate solutions. Let \mathbf{X}_\star be a solution to the model problem (2.2). For $\varepsilon \geq 0$, we say that a matrix \mathbf{X} is ε -optimal for (2.2) when

$$\|\mathcal{A}\mathbf{X} - \mathbf{b}\| \leq \varepsilon \quad \text{and} \quad \langle \mathbf{C}, \mathbf{X} \rangle - \langle \mathbf{C}, \mathbf{X}_\star \rangle \leq \varepsilon.$$

Many optimization algorithms aim to produce ε -optimal points.

As we saw in [subsection 1.3.2](#), there are many situations where the solutions to (2.2) have low rank, or they admit accurate low-rank approximations. The ε -optimal points inherit these properties for sufficiently small ε . This insight suggests a new computational goal.

For a rank parameter r , we will seek a rank- r matrix $\widehat{\mathbf{X}}$ that approximates an ε -optimal point \mathbf{X} . More precisely, for a fixed suboptimality parameter $\zeta > 0$, we want

$$(2.3) \quad \|\mathbf{X} - \widehat{\mathbf{X}}\|_* \leq (1 + \zeta) \cdot \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_* \quad \text{where} \quad \text{rank } \widehat{\mathbf{X}} \leq r \quad \text{and} \quad \mathbf{X} \text{ is } \varepsilon\text{-optimal}.$$

Given (2.3), if the ε -optimal point \mathbf{X} is close to *any* rank- r matrix, then the rank- r approximate solution $\widehat{\mathbf{X}}$ is also close to the ε -optimal point \mathbf{X} . This formulation is advantageous because it is easier to compute and to store the low-rank matrix $\widehat{\mathbf{X}}$.

2.3. Black-box presentation of problem data. To develop scalable algorithms for (2.2), it is productive to hide the internal complexity of the problem instance from the algorithm [115]. To do so, we treat \mathbf{C} and \mathcal{A} as black boxes that support three primitive computations:

$$(2.4) \quad \begin{array}{|c|} \hline \textcircled{1} \quad \mathbf{u} \mapsto \mathbf{C}\mathbf{u} \\ \mathbb{R}^n \rightarrow \mathbb{R}^n \\ \hline \end{array} \quad \begin{array}{|c|} \hline \textcircled{2} \quad (\mathbf{u}, \mathbf{z}) \mapsto (\mathcal{A}^* \mathbf{z})\mathbf{u} \\ \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}^n \\ \hline \end{array} \quad \begin{array}{|c|} \hline \textcircled{3} \quad \mathbf{u} \mapsto \mathcal{A}(\mathbf{u}\mathbf{u}^*) \\ \mathbb{R}^n \rightarrow \mathbb{R}^d \\ \hline \end{array}$$

The vectors $\mathbf{u} \in \mathbb{R}^n$ and $\mathbf{z} \in \mathbb{R}^d$ are arbitrary. Although these functions may seem abstract, they are often quite natural and easy to implement well. For example, see subsection 2.5.

We will formulate algorithms for (2.2) that interact with the problem data only through the operations (2.4). We tacitly assume that the primitives require minimal storage and arithmetic; otherwise, it may be impossible to develop a truly efficient algorithm.

Remark 2.1 (Inexact oracles). We assume the primitives (2.4) are implemented exactly. An interesting future direction is to study approximate or stochastic versions.

2.4. Resource usage. Our goal is to develop a scalable algorithm for the template (2.2) by computing a rank- r approximate solution that satisfies (2.3). Running time is less of an issue than storage because modern computers have fast processors but limited memory.

So what can we hope to achieve in terms of storage? First, observe that we need to specify the constraint vector $\mathbf{b} \in \mathbb{R}^d$ to determine an arbitrary instance of the model problem (2.2). Second, observe that a rank- r matrix with dimension n has $\Theta(rn)$ degrees of freedom. Therefore, we must be prepared to spend $\Theta(d + rn)$ numbers just to express the input and output of an SDP solver. The working storage of SketchyCGAL achieves this bound. This cost should be viewed in contrast to the $\Theta(n^2)$ degrees of freedom in the matrix variable in (2.2).

We tabulate the arithmetic costs of an algorithm by counting the number of times we apply each primitive, plus the total amount of extra arithmetic. For lower bounds on the arithmetic operations needed to solve an SDP to moderate accuracy, see [44, Thm. 2].

Remark 2.2 (Communication). Data transfer is the dominant expense in most large-scale applications. We do not control communication costs directly, but limiting the storage and arithmetic has the ancillary benefit of reducing data transfer.

2.5. Example: The MaxCut SDP. To provide a concrete example, let us summarize the meaning of the primitives and the desired storage costs for solving the MaxCut SDP (1.3).

- The primitive $\textcircled{1} : \mathbf{u} \mapsto -\mathbf{L}\mathbf{u}$. In the typical case that the Laplacian \mathbf{L} is sparse, this amounts to a sparse matrix–vector multiply.
- The primitive $\textcircled{2} : (\mathbf{u}, \mathbf{z}) \mapsto (\text{diag}^* \mathbf{z})\mathbf{u} = (z_1 u_1, \dots, z_n u_n)$.
- The primitive $\textcircled{3} : \mathbf{u} \mapsto \text{diag}(\mathbf{u}\mathbf{u}^*) = (|u_1|^2, \dots, |u_n|^2)$.
- The MaxCut SDP has n linear constraints, and we seek a rank-one approximation of the solution. Thus, we desire an algorithm that operates with $\Theta(n)$ working storage.

Note that we still need $\Theta(m)$ numbers to store the Laplacian \mathbf{L} of a generic graph with m edges. But we do not charge the optimization algorithm for this storage because the algorithm only interacts with \mathbf{L} through the primitives (2.4).

3. An algorithm for the model problem. We will develop a scalable method for the model problem (2.2) by enhancing an existing algorithm, called CGAL [110], developed by a subset of the authors. This method works well but lacks strong storage and arithmetic guarantees. This section summarizes the CGAL method and its convergence properties. Subsequent sections introduce additional ideas that we need to control resource usage, culminating with the SketchyCGAL algorithm in section 6.

3.1. The augmented problem. For a parameter $\beta > 0$, we revise the problem (2.2) by introducing an extra term in the objective:

$$(3.1) \quad \text{minimize} \quad \langle \mathbf{C}, \mathbf{X} \rangle + \frac{\beta}{2} \|\mathcal{A}\mathbf{X} - \mathbf{b}\|^2 \quad \text{subject to} \quad \mathcal{A}\mathbf{X} = \mathbf{b}, \quad \mathbf{X} \in \alpha\Delta_n.$$

The original problem (2.2) and the augmented problem (3.1) share the same optimal value and optimal set. But the new formulation has several benefits: the augmented objective cooperates with the affine constraint to penalize infeasible points, and the dual of the augmented problem is smoother than the dual of the original problem. See [20, Ch. 2] for background.

3.2. The augmented Lagrangian. We construct the Lagrangian L_β of the augmented problem (3.1) by introducing a dual variable $\mathbf{y} \in \mathbb{R}^d$ and promoting the affine constraint:

$$(3.2) \quad L_\beta(\mathbf{X}; \mathbf{y}) := \langle \mathbf{C}, \mathbf{X} \rangle + \langle \mathbf{y}, \mathcal{A}\mathbf{X} - \mathbf{b} \rangle + \frac{\beta}{2} \|\mathcal{A}\mathbf{X} - \mathbf{b}\|^2 \quad \text{for } \mathbf{X} \in \alpha\Delta_n \text{ and } \mathbf{y} \in \mathbb{R}^d.$$

For reference, note that the partial derivatives of the augmented Lagrangian L_β satisfy

$$(3.3) \quad \begin{aligned} \partial_{\mathbf{X}} L_\beta(\mathbf{X}; \mathbf{y}) &= \mathbf{C} + \mathcal{A}^* \mathbf{y} + \beta \mathcal{A}^* (\mathcal{A}\mathbf{X} - \mathbf{b}) \\ \partial_{\mathbf{y}} L_\beta(\mathbf{X}; \mathbf{y}) &= \mathcal{A}\mathbf{X} - \mathbf{b}. \end{aligned}$$

We attempt to minimize the augmented Lagrangian L_β with respect to the primal variable \mathbf{X} , while we attempt to maximize with respect to the dual variable \mathbf{y} .

3.3. The augmented Lagrangian strategy. The form of the augmented Lagrangian suggests an algorithm. We generate a sequence $\{(\mathbf{X}_t; \mathbf{y}_t)\}$ of primal–dual pairs by alternately minimizing over the primal variable \mathbf{X} and taking a gradient step in the dual variable \mathbf{y} .

1. **Initialize:** Choose $\mathbf{X}_0 \in \alpha\Delta_n$ and $\mathbf{y}_0 \in \mathbb{R}^d$.
2. **Primal step:** $\mathbf{X}_{t+1} \in \arg \min \{L_\beta(\mathbf{X}; \mathbf{y}_t) : \mathbf{X} \in \alpha\Delta_n\}$.
3. **Dual step:** $\mathbf{y}_{t+1} = \mathbf{y}_t + \beta(\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b})$.

As we proceed, we can also increase the smoothing parameter β to make violations of the affine constraint in (3.1) more and more intolerable.

The augmented Lagrangian strategy is powerful, but it is hard to apply in this setting because of the cost of implementing the primal step, even approximately.

3.4. The CGAL iteration. The CGAL iteration [110] is related to the augmented Lagrangian (AL) paradigm, but the primal steps are inspired by the conditional gradient method (CGM). CGAL identifies an update direction for the primal variable by minimizing a linear proxy for the augmented Lagrangian. We take a small primal step in the update direction,

and we improve the dual variable by taking a small gradient step. At each iteration, the smoothing parameter increases according to a fixed schedule.

This subsection outlines the steps in the CGAL iteration. Afterward, we give *a priori* guarantees on the convergence rate. Pseudocode for CGAL appears in [Algorithm A.1](#).

3.4.1. Initialization. Let $\beta_0 > 0$ be an initial smoothing parameter, and fix a (large) bound $K > 0$ on the maximum allowable size of the dual variable. We also assume that we have access to the norm $\|\mathcal{A}\|$ of the constraint matrix, or—failing that—a *lower* bound.

Begin with an arbitrary choice $\mathbf{X}_1 \in \mathbb{S}_n$ for the primal matrix variable, and set the initial dual vector variable $\mathbf{y}_1 = \mathbf{0}$.

3.4.2. Primal updates via linear minimization. At iteration $t = 1, 2, 3, \dots$, we increase the smoothing parameter to $\beta_t := \beta_0 \sqrt{t+1}$, and we form the partial derivative of the augmented Lagrangian with respect to the primal variable at the current pair of iterates:

$$(3.4) \quad \mathbf{D}_t := \partial_{\mathbf{X}} L_{\beta_t}(\mathbf{X}_t; \mathbf{y}_t) = \mathbf{C} + \mathcal{A}^*(\mathbf{y}_t + \beta_t(\mathcal{A}\mathbf{X}_t - \mathbf{b})).$$

Then we compute an update $\mathbf{H}_t \in \alpha \Delta_t$ by finding the feasible point that is most correlated with the negative gradient $-\mathbf{D}_t$. This amounts to a linear minimization problem:

$$(3.5) \quad \begin{aligned} \mathbf{H}_t &\in \arg \min \{ \langle \mathbf{D}_t, \mathbf{H} \rangle : \mathbf{H} \in \alpha \Delta_n \} \\ &= \text{convex hull} \{ \alpha \mathbf{v} \mathbf{v}^* : \mathbf{v} \text{ is a unit-norm minimum eigenvector of } \mathbf{D}_t \}. \end{aligned}$$

In particular, we may take $\mathbf{H}_t = \alpha \mathbf{v}_t \mathbf{v}_t^*$ for a minimum eigenvector \mathbf{v}_t of the gradient \mathbf{D}_t . Next, update the matrix primal variable by taking a small step in the direction \mathbf{H}_t :

$$(3.6) \quad \mathbf{X}_{t+1} := (1 - \eta_t) \mathbf{X}_t + \eta_t \mathbf{H}_t \in \alpha \Delta_n \quad \text{where} \quad \eta_t := \frac{2}{t+1}.$$

The appeal of this approach is that it only requires a single eigenvector computation (3.5). Moreover, we need not compute the eigenvector accurately; see [subsection 3.4.4](#) for details.

3.4.3. Dual updates via gradient ascent. Next, we update the dual variable by taking a small gradient step on the dual variable in the augmented Lagrangian:

$$(3.7) \quad \mathbf{y}_{t+1} = \mathbf{y}_t + \gamma_t \partial_{\mathbf{y}} L_{\beta_t}(\mathbf{X}_{t+1}; \mathbf{y}_t) = \mathbf{y}_t + \gamma_t (\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}).$$

The dual step size γ_t is the largest number that satisfies the conditions

$$(3.8) \quad \gamma_t \|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\|^2 \leq \frac{4\alpha^2\beta_0}{(t+1)^{3/2}} \|\mathcal{A}\|^2 \quad \text{and} \quad 0 \leq \gamma_t \leq \beta_0.$$

We omit the dual step if it makes the dual variable too large. More precisely, if (3.7) and (3.8) result in $\|\mathbf{y}_{t+1}\| > K$, then we set $\mathbf{y}_{t+1} = \mathbf{y}_t$ instead. This is the CGAL iteration.

3.4.4. Approximate eigenvector computations. It is crucial that the CGAL strategy provably works if we replace (3.5) with an approximate minimum eigenvector computation. At step t , suppose that we find an update direction

$$(3.9) \quad \mathbf{H}_t := \alpha \mathbf{v}_t \mathbf{v}_t^* \quad \text{where} \quad \mathbf{v}_t^* \mathbf{D}_t \mathbf{v}_t \leq \lambda_{\min}(\mathbf{D}_t) + \frac{1}{\sqrt{t+1}} \|\mathbf{D}_t\|$$

for a unit vector \mathbf{v}_t . As usual, λ_{\min} is the minimum eigenvalue. The matrix \mathbf{H}_t from (3.9) serves in place of a solution to (3.5). We explain how to solve the subproblem (3.9) in [section 4](#).

3.4.5. Solution quality. Given a primal–dual pair $(\mathbf{X}_t; \mathbf{y}_t)$, we can assess the quality of the primal solution to the model problem (2.2) and decide when to halt. See Appendix A.4.

3.5. Convergence guarantees for CGAL. The following result describes the convergence of the CGAL iteration. The analysis is adapted from [110, Thm. 3.1]; see Appendix A for details and a recapitulation of the proof.

Fact 3.1 (CGAL: Convergence). *Assume problem (2.2) satisfies strong duality. The CGAL iteration (subsection 3.4) with approximate eigenvector computations (3.9) yields a sequence $\{\mathbf{X}_t : t = 1, 2, 3, \dots\} \subset \alpha \Delta_n$ that satisfies*

$$(3.10) \quad \|\mathcal{A}\mathbf{X}_t - \mathbf{b}\| \leq \frac{\text{Const}}{\sqrt{t}} \quad \text{and} \quad |\langle \mathbf{C}, \mathbf{X}_t \rangle - \langle \mathbf{C}, \mathbf{X}_\star \rangle| \leq \frac{\text{Const}}{\sqrt{t}}.$$

The matrix \mathbf{X}_\star solves (2.2). The constant depends on the problem data $(\mathbf{C}, \mathcal{A}, \mathbf{b}, \alpha)$, the minimum Euclidean norm of a dual solution, and the algorithm parameters β_0 and K ,

In view of (3.10), CGAL produces a primal iterate \mathbf{X}_T that is ε -optimal after $T = O(\varepsilon^{-2})$ iterations. The numerical work in [110, Sec. 5] shows that CGAL usually achieves an ε -optimal point after $T = O(\varepsilon^{-1})$ iterations.

The analysis behind Fact 3.1 indicates that CGAL converges more quickly if we precondition the problem data by rescaling; see subsection 7.1.2. This step is critical in practice.

3.6. Distributed computation. Since CGAL builds on the augmented Lagrangian framework, we can apply it even when the problem is too large to solve on one computational node. In particular, when d is large, it may be advantageous to partition the constraint matrices \mathbf{A}_i and the associated dual variables y_i among several workers. Distributed CGAL has a structure similar to the alternating directions method of multipliers (ADMM) [24].

4. Approximate eigenvectors. Most of the computation in CGAL occurs in the approximate eigenvector step (3.9). Here, we encounter a tradeoff between arithmetic and storage.

4.1. Krylov methods. The approximate minimum eigenvector computation (3.9) involves a matrix of the form $\mathbf{D}_t = \mathbf{C} + \mathcal{A}^* \mathbf{w}_t$. Observe that we can multiply the matrix \mathbf{D}_t by any vector using the primitives (2.4) 1–2. Krylov methods compute eigenvectors of \mathbf{D}_t by repeated matrix–vector multiplication with \mathbf{D}_t , so they are obvious tools for the subproblem (3.9).

Unfortunately, CGAL tends to generate matrices \mathbf{D}_t that have clustered eigenvalues. These matrices challenge standard eigenvector software, such as ARPACK [70], the engine behind the MATLAB command `eigs`. Instead, we retreat to more classical techniques.

4.2. A randomized shifted power method. First, let us consider the situation where storage is the overriding concern. In this case, we run the shifted power method with a random starting vector; see Algorithm 4.1. Kuczyński & Woźniakowski [65, Thm. 4.1(a)] have obtained error bounds for this type of algorithm.

Fact 4.1 (Randomized shifted power method). *Let $\mathbf{M} \in \mathbb{S}_n$. For $\varepsilon \in (0, 1]$ and $\delta \in (0, 1]$, the shifted power method, Algorithm 4.1, computes a unit vector $\mathbf{u} \in \mathbb{F}^n$ that satisfies*

$$\mathbf{u}^* \mathbf{M} \mathbf{u} \leq \lambda_{\min}(\mathbf{M}) + \frac{3\varepsilon}{2} \|\mathbf{M}\| \quad \text{with probability at least } 1 - \delta$$

Algorithm 4.1 ApproxMinEvec via randomized shifted power method (subsection 4.2).

Input: Input matrix $M \in \mathbb{S}_n$ and maximum number q of iterations

Output: Approximate minimum eigenpair $(\xi, v) \in \mathbb{R} \times \mathbb{F}^n$ of the matrix M

```

1 function ApproxMinEvec( $M$ ;  $q$ )
2    $\sigma \leftarrow 2 \|M\|$  ▷ Use normest to estimate shift
3    $v \leftarrow \text{randn}(n, 1)/\sqrt{n}$  ▷ Random initial vector
4   for  $i \leftarrow 1, 2, 3, \dots, q$  do
5      $v \leftarrow \sigma v - Mv$  ▷ Power method for  $\sigma I - M$ 
6      $v \leftarrow v/\|v\|$  ▷ Approx. minimum eigenvector of  $M$ 
7    $\xi \leftarrow v^*(Mv)$  ▷ Approx. minimum eigenvalue of  $M$ 

```

after $q \geq \frac{1}{2} + \varepsilon^{-1} \log(n/\delta^2)$ iterations.¹ The arithmetic cost is $O(q)$ matrix–vector multiplies with M and $O(qn)$ extra operations. The working storage is about $2n$ numbers.

With constant probability, we solve the eigenvector problem (3.9) successfully in every iteration t of CGAL if we use $q_t = O(t^{1/2} \log(tn))$ iterations of Algorithm 4.1. In practice, we implement CGAL with the concrete choice $q_t = 8t^{1/2} \log n$, which behaves well empirically.

4.3. A randomized Lanczos method. Suppose we are willing to allocate a moderate amount of storage to perform the approximate eigenvector computation. In this case, we can run the Lanczos iteration with a random starting vector [48, Sec. 10.1]. See Algorithm 4.2 for pseudocode. Kuczyński & Woźniakowski [65, Thm. 4.2(a)] have obtained error bounds.

Fact 4.2 (Randomized Lanczos method). Let $M \in \mathbb{S}_n$. For $\varepsilon \in (0, 1]$ and $\delta \in (0, 0.5]$, the randomized Lanczos method, Algorithm 4.2, computes a unit vector $u \in \mathbb{F}^n$ that satisfies

$$u^* M u \leq \lambda_{\min}(M) + \frac{\varepsilon}{8} \|M\| \quad \text{with probability at least } 1 - 2\delta$$

after $q \geq \frac{1}{2} + \varepsilon^{-1/2} \log(n/\delta^2)$ iterations. The arithmetic cost is at most q matrix–vector multiplies with M and $O(qn + q^2)$ extra operations. The working storage is $O(qn)$.

With constant probability, we solve the eigenvector problem (3.9) successfully in every iteration t of CGAL if we use $q_t = O(t^{1/4} \log(tn))$ iterations of Algorithm 4.2. In practice, we implement CGAL with the explicit choice $q_t = t^{1/4} \log n$.

5. Sketching and reconstruction of a psd matrix. The CGAL iteration generates a psd matrix that solves the model problem (2.2) via a sequence (3.6) of rank-one linear updates. To control storage costs, SketchyCGAL will retain only a summary of the psd matrix variable. This section outlines the *Nystrom sketch*, an established method [49, 46, 72, 100] that can track the evolving psd matrix and then report a provably accurate low-rank approximation.

5.1. Sketching and updates. Consider a psd input matrix $X \in \mathbb{S}_n$. Let R be a parameter that modulates the storage cost of the sketch and the quality of the matrix approximation.

¹All logarithms are base-e.

Algorithm 4.2 ApproxMinEvec via randomized Lanczos method (subsection 4.3).

Input: Input matrix $M \in \mathbb{S}_n$ and maximum number q of iterations

Output: Approximate minimum eigenpair $(\xi, \mathbf{v}) \in \mathbb{R} \times \mathbb{F}^n$ of the matrix M

```

1 function ApproxMinEvec( $M; q$ )
2    $\mathbf{v}_1 \leftarrow \text{randn}(n, 1)$  ▷ Aggregate  $\mathbf{v}_i$  as columns of  $V$ 
3    $\mathbf{v}_1 \leftarrow \mathbf{v}_1 / \|\mathbf{v}_1\|$ 
4   for  $i \leftarrow 1, 2, 3, \dots, \min\{q, n-1\}$  do
5      $\alpha_i \leftarrow \text{Re}(\mathbf{v}_i^* (M \mathbf{v}_i))$ 
6      $\mathbf{v}_{i+1} \leftarrow M \mathbf{v}_i - \alpha_i \mathbf{v}_i - \beta_{i-1} \mathbf{v}_{i-1}$  ▷ Lanczos iteration;  $\beta_0 \mathbf{v}_0 = \mathbf{0}$ 
7     if  $\beta_i = 0$  then break ▷ Found an invariant subspace!
8      $\mathbf{v}_{i+1} \leftarrow \mathbf{v}_{i+1} / \beta_i$ 
9    $T \leftarrow \text{tridiag}(\beta_{1:(i-1)}, \alpha_{1:i}, \beta_{1:(i-1)})$  ▷ Form tridiagonal matrix
10   $[U, D] \leftarrow \text{eig}(T)$ 
11   $[\xi, \text{ind}] \leftarrow \min(\text{diag}(D))$  ▷ Approx. minimum eigenvalue
12   $\mathbf{v} \leftarrow V_{1:i} U(:, \text{ind}(1))$  ▷ Approx. minimum eigenvector

```

To construct the Nyström sketch, we draw and fix a standard normal² test matrix $\Omega \in \mathbb{F}^{n \times R}$. Our summary, or *sketch*, of the matrix X takes the form

$$(5.1) \quad S = X\Omega \in \mathbb{F}^{n \times R}.$$

The sketch supports linear rank-one updates to X . Indeed, we can track the evolution

$$(5.2) \quad \begin{aligned} X &\leftarrow (1 - \eta) X + \eta \mathbf{v} \mathbf{v}^* && \text{for } \eta \in [0, 1] \text{ and } \mathbf{v} \in \mathbb{F}^n \\ \text{via } S &\leftarrow (1 - \eta) S + \eta \mathbf{v} (\mathbf{v}^* \Omega). \end{aligned}$$

The test matrix Ω and the sketch S require storage of $2Rn$ numbers in \mathbb{F} . The arithmetic cost of the linear update (5.2) to the sketch is $\Theta(Rn)$ numerical operations.

Remark 5.1 (Structured random matrices). We can reduce storage costs by a factor of two by using a structured random matrix in place of Ω . For example, see [102, Sec. 3] or [98].

5.2. The reconstruction process. Given the test matrix Ω and the sketch $S = X\Omega$, we form a rank- R approximation \widehat{X} of the sketched matrix X . The approximation is defined by

$$(5.3) \quad \widehat{X} := S(\Omega^* S)^\dagger S^* = (X\Omega)(\Omega^* X\Omega)^\dagger (X\Omega)^*,$$

where † is the pseudoinverse. This reconstruction is called a *Nyström approximation*. We often truncate \widehat{X} by replacing it with its best rank- r approximation $[\widehat{X}]_r$ for some $r \leq R$.

See Algorithm 5.1 for a numerically stable implementation of the Nyström approximation (5.3). The algorithm takes $\Theta(R^2 n)$ numerical operations and $\Theta(Rn)$ storage.

²Each entry of the matrix is an independent Gaussian random variable with mean zero and variance one. In the complex setting, the real and imaginary parts of each entry are independent standard normal variables.

Algorithm 5.1 NystromSketch implementation (see [section 5](#))**Input:** Dimension n of input matrix, size R of sketch**Output:** Rank- R approximation $\widehat{\mathbf{X}}$ of sketched matrix in factored form $\widehat{\mathbf{X}} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$, where $\mathbf{U} \in \mathbb{F}^{n \times R}$ has orthonormal columns and $\mathbf{\Lambda} \in \mathbb{R}^{R \times R}$ is nonnegative diagonal**Recommendation:** Choose R as large as possible, given storage and arithmetic constraints

```

1 function NystromSketch.Init( $n, R$ )
2    $\mathbf{\Omega} \leftarrow \text{randn}(n, R)$  ▷ Draw and fix random test matrix
3    $\mathbf{S} \leftarrow \text{zeros}(n, R)$  ▷ Form sketch of zero matrix

4 function NystromSketch.RankOneUpdate( $\mathbf{v}, \eta$ ) ▷ Implements (5.2)
5    $\mathbf{S} \leftarrow (1 - \eta) \mathbf{S} + \eta \mathbf{v}(\mathbf{v}^* \mathbf{\Omega})$  ▷ Update sketch of matrix

6 function NystromSketch.Reconstruct()
7    $\sigma \leftarrow \sqrt{n} \text{eps}(\text{norm}(\mathbf{S}))$  ▷ Compute a shift parameter
8    $\mathbf{S}_\sigma \leftarrow \mathbf{S} + \sigma \mathbf{\Omega}$  ▷ Implicitly form sketch of  $\mathbf{X} + \sigma \mathbf{I}$ 
9    $\mathbf{C} \leftarrow \text{chol}(\mathbf{\Omega}^* \mathbf{S}_\sigma)$ 
10   $(\mathbf{U}, \mathbf{\Sigma}, \sim) \leftarrow \text{svd}(\mathbf{S}_\sigma / \mathbf{C})$  ▷ Dense SVD
11   $\mathbf{\Lambda} \leftarrow \max\{0, \mathbf{\Sigma}^2 - \sigma \mathbf{I}\}$  ▷ Remove shift

```

5.3. A priori error bounds. The Nyström approximation $\widehat{\mathbf{X}}$ yields a provably good estimate for the matrix \mathbf{X} contained in the sketch [100, Thm. 4.1].

Fact 5.2 (Nyström sketch: Error bound). Fix a psd matrix $\mathbf{X} \in \mathbb{S}_n$. Let $\mathbf{S} = \mathbf{X}\mathbf{\Omega}$ where $\mathbf{\Omega} \in \mathbb{F}^{n \times R}$ is standard normal. For each $r < R$, the Nyström approximation (5.3) satisfies

$$(5.4) \quad \mathbb{E}_{\mathbf{\Omega}} \|\mathbf{X} - \widehat{\mathbf{X}}\|_* \leq \left(1 + \frac{r}{R - r - 1}\right) \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_*,$$

where $\mathbb{E}_{\mathbf{\Omega}}$ is the expectation with respect to $\mathbf{\Omega}$. If we replace $\widehat{\mathbf{X}}$ with its rank- r truncation $\llbracket \widehat{\mathbf{X}} \rrbracket_r$, the error bound (5.4) remains valid. Similar results hold with high probability.

There is a lot more to say about the implementation and behavior of low-rank matrix approximation from streaming data. See [Appendix B](#) and our papers [100, 101, 102].

6. Scalable semidefinite programming via SketchyCGAL. We may now present an extension of CGAL that solves the model problem (2.2) with controlled storage and computation.

6.1. The opportunity. Our new algorithm SketchyCGAL augments the CGAL iteration from [subsection 3.4](#). Instead of storing the matrix \mathbf{X}_t in the CGAL iteration...

1. We drive the iteration with the d -dimensional primal state variable $\mathbf{z}_t := \mathcal{A}\mathbf{X}_t$.
2. We maintain a Nyström sketch of the primal iterate \mathbf{X}_t using storage $\Theta(Rn)$.
3. When the iteration is halted, say at step T , we use the sketch to construct a rank- R approximation $\widehat{\mathbf{X}}_T$ of the implicitly computed solution \mathbf{X}_T of the model problem.

As we will see, the resulting method exhibits almost the same convergence behavior as the CGAL algorithm, but it also enjoys a strong storage guarantee.

6.2. The SketchyCGAL iteration. To develop the SketchyCGAL iteration, we begin with the CGAL iteration. Then make the substitutions $\mathbf{z}_t = \mathcal{A}\mathbf{X}_t$ and $\mathbf{h}_t = \mathcal{A}\mathbf{H}_t$ to eliminate the matrix variables. Let us summarize what happens; see subsection 6.3 for additional explanation. Algorithm 6.1 contains pseudocode with implementation recommendations.

6.2.1. Initialization. First, we choose the size R of the Nyström sketch. Then draw and fix the random test matrix $\mathbf{\Omega} \in \mathbb{F}^{n \times R}$. Select an initial smoothing parameter β_0 and a bound K on the dual variable. Define the smoothing and step size parameters:

$$(6.1) \quad \beta_t := \beta_0 \sqrt{t+1} \quad \text{and} \quad \eta_t := \frac{2}{t+1}.$$

Initialize the primal state variable, the sketch, and the dual variable:

$$(6.2) \quad \mathbf{z}_1 := \mathbf{0} \in \mathbb{R}^d \quad \text{and} \quad \mathbf{S}_1 := \mathbf{0} \in \mathbb{F}^{n \times R} \quad \text{and} \quad \mathbf{y}_1 := \mathbf{0} \in \mathbb{R}^d.$$

These choices correspond to the simplest initial iterate $\mathbf{X}_1 := \mathbf{0}$.

6.2.2. Primal updates. At iteration $t = 1, 2, 3, \dots$, we compute a unit-norm vector \mathbf{v}_t that is an approximate minimum eigenvector of the gradient \mathbf{D}_t of the smoothed objective:

$$(6.3) \quad \mathbf{v}_t^* \mathbf{D}_t \mathbf{v}_t \leq \lambda_{\min}(\mathbf{D}_t) + \frac{\|\mathbf{D}_t\|}{\sqrt{t+1}} \quad \text{where} \quad \mathbf{D}_t := \mathbf{C} + \mathcal{A}^*(\mathbf{y}_t + \beta_t(\mathbf{z}_t - \mathbf{b})).$$

This calculation corresponds with (3.9). Form the primal update direction $\mathbf{h}_t = \mathcal{A}(\alpha \mathbf{v}_t \mathbf{v}_t^*)$, and then update the primal state variable \mathbf{z}_t and the sketch \mathbf{S}_t :

$$(6.4) \quad \mathbf{z}_{t+1} := (1 - \eta_t)\mathbf{z}_t + \eta_t \mathbf{h}_t \quad \text{and} \quad \mathbf{S}_{t+1} := (1 - \eta_t)\mathbf{S}_t + \eta_t \alpha \mathbf{v}_t (\mathbf{v}_t^* \mathbf{\Omega}).$$

We obtain the update rule for the primal state variable \mathbf{z}_t by applying the linear map \mathcal{A} to the primal update rule (3.6). The update rule for the sketch \mathbf{S}_t follows from (5.2).

6.2.3. Dual updates. The update to the dual variable takes the form

$$(6.5) \quad \mathbf{y}_{t+1} = \mathbf{y}_t + \gamma_t(\mathbf{z}_{t+1} - \mathbf{b})$$

where we choose the largest γ_t that satisfies the conditions

$$(6.6) \quad \gamma_t \|\mathbf{z}_{t+1} - \mathbf{b}\|^2 \leq \frac{4\alpha^2 \beta_0}{(t+1)^{3/2}} \|\mathcal{A}\|^2 \quad \text{and} \quad 0 \leq \gamma_t \leq \beta_0.$$

If needed, we set $\gamma_t = 0$ to prevent $\|\mathbf{y}_{t+1}\| > K$. This is the SketchyCGAL iteration.

6.3. Connection with CGAL. There is a tight connection between the iterates of SketchyCGAL and CGAL. Let $\mathbf{X}_1 := \mathbf{0}$. Using the vectors \mathbf{v}_t computed in (6.3), define matrices

$$(6.7) \quad \mathbf{H}_t := \alpha \mathbf{v}_t \mathbf{v}_t^* \quad \text{and} \quad \mathbf{X}_{t+1} := (1 - \eta_t)\mathbf{X}_t + \eta_t \mathbf{H}_t.$$

With these definitions, the following loop invariants are in force:

$$(6.8) \quad \mathbf{h}_t = \mathcal{A}\mathbf{H}_t \quad \text{and} \quad \mathbf{z}_t = \mathcal{A}\mathbf{X}_t \quad \text{and} \quad \mathbf{S}_t = \mathbf{X}_t \mathbf{\Omega}.$$

By comparing subsections 3.4 and 6.2, we see that the trajectory $\{(\mathbf{X}_t, \mathbf{H}_t, \mathbf{y}_t) : t = 1, 2, 3, \dots\}$ could also have been generated by running the CGAL iteration. In other words, the variables in SketchyCGAL track the variables of some invocation of CGAL and inherit their behavior. We refer to the matrices \mathbf{X}_t as the *implicit* CGAL iterates.

Algorithm 6.1 SketchyCGAL for the model problem (2.2)

Input: Problem data for (2.2) implemented via the primitives (2.4), sketch size R , number T of iterations

Output: Rank- R approximate solution to (2.2) in factored form $\widehat{\mathbf{X}}_T = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$ where $\mathbf{U} \in \mathbb{R}^{n \times R}$ has orthonormal columns and $\mathbf{\Lambda} \in \mathbb{R}^{R \times R}$ is nonnegative diagonal

Recommendation: To achieve (2.3), set R as large as possible, and set $T \approx \varepsilon^{-1}$

```

1 function SketchyCGAL( $R; T$ )
2   Scale problem data (subsection 7.1.2)                                ▷ [opt] Recommended!
3    $\beta_0 \leftarrow 1$  and  $K \leftarrow +\infty$                                ▷ Default parameters
4   NystromSketch.Init( $n, R$ )
5    $\mathbf{z} \leftarrow \mathbf{0}_d$  and  $\mathbf{y} \leftarrow \mathbf{0}_d$ 
6   for  $t \leftarrow 1, 2, 3, \dots, T$  do
7      $\beta \leftarrow \beta_0 \sqrt{t+1}$  and  $\eta \leftarrow 2/(t+1)$ 
8      $(\xi, \mathbf{v}) \leftarrow \text{ApproxMinEvec}(\mathbf{C} + \mathcal{A}^*(\mathbf{y} + \beta(\mathbf{z} - \mathbf{b})); q_t)$     ▷ Use primitives (2.4) ①②
                                           ▷ With Algorithm 4.1, set  $q_t = 8t^{1/2} \log n$ .
                                           ▷ With Algorithm 4.2, set  $q_t = t^{1/4} \log n$ 
9      $\mathbf{z} \leftarrow (1 - \eta) \mathbf{z} + \eta \mathcal{A}(\alpha \mathbf{v} \mathbf{v}^*)$                                 ▷ Use primitive (2.4) ③
10     $\mathbf{y} \leftarrow \mathbf{y} + \gamma(\mathbf{z} - \mathbf{b})$                                            ▷  $\gamma$  is the largest solution to (6.6)
11    NystromSketch.RankOneUpdate( $\sqrt{\alpha} \mathbf{v}, \eta$ )
12     $(\mathbf{U}, \mathbf{\Lambda}) \leftarrow \text{NystromSketch.Reconstruct}$ 
13     $\mathbf{\Lambda} \leftarrow \mathbf{\Lambda} + (\alpha - \text{tr}(\mathbf{\Lambda})) \mathbf{I}_R / R$                     ▷ [opt] Enforce trace constraint in (2.2)
```

6.3.1. Approximating the CGAL iterates. We do not have access to the implicit CGAL iterates described above. Nevertheless, the sketch permits us to approximate them! After iteration t of SketchyCGAL, we can form a rank- R approximation $\widehat{\mathbf{X}}_t$ of the implicit iterate \mathbf{X}_t by invoking the formula (5.3) with $\mathbf{S} = \mathbf{S}_t$. According to Fact 5.2, for each $r < R$,

$$(6.9) \quad \mathbb{E}_\Omega \|\mathbf{X}_t - \widehat{\mathbf{X}}_t\|_* \leq \left(1 + \frac{r}{R - r - 1}\right) \|\mathbf{X}_t - \llbracket \mathbf{X}_t \rrbracket_r\|_*.$$

In other words, the computed approximation $\widehat{\mathbf{X}}_t$ is a good proxy for the implicit iterate \mathbf{X}_t whenever the latter matrix is well-approximated by a low-rank matrix. The same bound (6.9) holds if we replace $\widehat{\mathbf{X}}_t$ by the truncated rank- r matrix $\llbracket \widehat{\mathbf{X}}_t \rrbracket_r$.

Remark 6.1 (Trace constraint). The model problem (2.2) requires the trace of the matrix variable to equal α , but the computed solution $\widehat{\mathbf{X}}_t$ rarely satisfies this constraint. Algorithm 6.1 includes an optional step that corrects the trace of the computed solution. Although this step is principled and justifiable, our theoretical analysis *does not* include its effects.

6.3.2. Solution quality. Given the quantities computed by SketchyCGAL, we can assess how well the implicit iterate \mathbf{X}_t solves the model problem (2.2). See Appendix C.1.

6.4. Convergence of SketchyCGAL. The implicit iterates \mathbf{X}_T converge to a solution of the model problem (2.2) at the same rate as the iterates of CGAL would. On average, the

rank- R iterates $\widehat{\mathbf{X}}_T$ track the implicit iterates. The discrepancy between them depends on how well the implicit iterates are approximated by low-rank matrices. Here is a simple convergence result that reflects this intuition; see [Appendix C.2](#) for the proof and further results.

Theorem 6.2 (SketchyCGAL: Convergence). *Assume problem (2.2) satisfies strong duality, and let Ψ_* be its solution set. For each $r < R$, the iterates $\widehat{\mathbf{X}}_t$ computed by SketchyCGAL (subsections 6.2 and 6.3) satisfy*

$$\limsup_{t \rightarrow \infty} \mathbb{E}_{\Omega} \text{dist}_*(\widehat{\mathbf{X}}_t, \Psi_*) \leq \left(1 + \frac{r}{R - r - 1}\right) \cdot \max_{\mathbf{X} \in \Psi_*} \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_*.$$

The same bound holds if we replace $\widehat{\mathbf{X}}_t$ with the truncated approximation $\llbracket \widehat{\mathbf{X}}_t \rrbracket_r$. Here, dist_* is the nuclear-norm distance between a matrix and a set of matrices.

6.5. Resource usage. We may now account for the resources required by SketchyCGAL.

6.5.1. Working storage. The main working variables in the SketchyCGAL iteration are the primal state $\mathbf{z}_t \in \mathbb{R}^d$, the computed eigenvector $\mathbf{v}_t \in \mathbb{R}^n$, the update direction $\mathbf{h}_t \in \mathbb{R}^d$, the test matrix $\Omega \in \mathbb{R}^{n \times R}$, the sketch $\mathbf{S}_t \in \mathbb{R}^{n \times R}$, and the dual variable $\mathbf{y}_t \in \mathbb{R}^d$. The total cost for storing these variables is $\Theta(d + Rn)$.

6.5.2. Approximate eigenvectors. The bottleneck in SketchyCGAL is the approximate eigenvector computation (6.3). For the subproblem in iteration t , we have two options:

1. **Shifted power.** Perform $q_t = O(t^{1/2} \log(tn))$ steps of [Algorithm 4.1](#), expending $O(q_t)$ applications of primitives (2.4) 1 2, plus $O(q_t n)$ arithmetic and $O(n)$ storage.
2. **Randomized Lanczos.** Do $q_t = O(t^{1/4} \log(tn))$ steps of [Algorithm 4.2](#), using $O(q_t)$ applications of primitives (2.4) 1 2, plus $O(q_t n + q_t^2)$ arithmetic and $O(q_t n)$ storage.

These choices ensure that, with constant probability, the eigenvector computation succeeds in every iteration t of SketchyCGAL. In a practical implementation, we undertake somewhat fewer steps in the eigenvector algorithms; see subsections 4.2 and 4.3.

6.5.3. Variable updates. The remaining computation takes place in the variable updates. To form the primal update direction \mathbf{h}_t , we invoke the primitive (2.4) 3. To update the primal state variable \mathbf{z}_t and the sketch \mathbf{S}_t in (6.4), we need $O(d + Rn)$ arithmetic operations. No further storage is required at this stage.

6.5.4. Overview of costs. Table 1 documents the cost of performing T iterations of the SketchyCGAL method. The first column summarizes the resources consumed in the outer iteration. The second column lists the additional resource usage if we invoke the shifted power method ([Algorithm 4.1](#)) to solve the approximate eigenvalue problem (6.3). The third column lists the extra resources if instead we invoke the Lanczos method ([Algorithm 4.2](#)).

In light of [Fact 3.1](#) and [subsection 6.3](#), we can be confident that the implicit iterate \mathbf{X}_T is ε -optimal within $T = O(\varepsilon^{-2})$ iterations. The formula (6.9) gives *a priori* guarantees on the quality of the approximation $\widehat{\mathbf{X}}_T$.

6.6. Theoretical performance of SketchyCGAL. We can package up this discussion in a theorem that describes the theoretical performance of the SketchyCGAL method.

Table 1

Resource usage for T iterations of *SketchyCGAL* with sketch size R using two different eigensolvers. The algorithm returns a rank- R approximation to an ε -optimal solution (subsection 2.2) to the model problem (2.2). In theory, $T \sim \varepsilon^{-2}$. In practice, $T \sim \varepsilon^{-1}$. Constants and logarithms are omitted. See subsection 6.5.4.

	SketchyCGAL Algorithm 6.1	w/shifted power Algorithm 4.1	or w/Lanczos Algorithm 4.2
Storage (floats)	$d + Rn$	n	$T^{1/4}n$
Arithmetic (flops)	$T(d + Rn)$	$T^{3/2}n$	$T^{5/4}n + T^{3/2}$
Primitives (calls)			
(2.4) 1 2	—	$T^{3/2}$	$T^{5/4}$
(2.4) 3	T	—	—

Theorem 6.3 (SketchyCGAL). Assume the model problem (2.2) satisfies strong duality. Fix a rank r , and set the sketch size $R \geq \zeta^{-1}r + 1$. After $T = O(\varepsilon^{-2})$ iterations, *SketchyCGAL* (subsections 6.2 and 6.3) returns a rank- r approximation $\widehat{\mathbf{X}} = \llbracket \widehat{\mathbf{X}}_T \rrbracket_r$ to an ε -optimal point of (2.2) that satisfies (2.3) with constant probability.

1. **Shifted power.** Suppose we use Algorithm 4.1 for the subproblem (6.3). The storage cost is $O(d + Rn)$ floats. Arithmetic costs are $O(\varepsilon^{-2}(d + Rn) + \varepsilon^{-3}n \log(n/\varepsilon))$ flops, $O(\varepsilon^{-3} \log(n/\varepsilon))$ invocations of primitives (2.4) 1 2, and $O(\varepsilon^{-2})$ applications of (2.4) 3.
2. **Lanczos.** Suppose we use Algorithm 4.2 to solve (6.3). Then the storage cost is $O(d + Rn + \varepsilon^{-1/2}n \log(n/\varepsilon))$ floats. The arithmetic requirements are $O(\varepsilon^{-2}(d + Rn) + (\varepsilon^{-3/2}n + \varepsilon^{-3}) \log(n/\varepsilon))$ flops, $O(\varepsilon^{-5/2} \log(n/\varepsilon))$ invocations of the primitives (2.4) 1 2, and $O(\varepsilon^{-2})$ applications of the primitive (2.4) 3.

In practice, *SketchyCGAL* performs substantially better than Theorem 6.3 suggests: empirically, we find that $T = O(\varepsilon^{-1})$. See section 7 for numerical evidence.

6.7. Comparison with CGAL. Like CGAL, *SketchyCGAL* is a general SDP solver with strong guarantees. Yet our changes lead to a new method with a better computational profile.

We motivated sketching as a mechanism to control storage. This improvement is visible when we compare the $\Theta(n^2)$ storage cost of CGAL to the $\Theta(d + Rn)$ storage cost of *SketchyCGAL*. But sketching also reduces arithmetic and communication. Indeed, *SketchyCGAL* works with state vectors of length d , and it updates the $\Theta(Rn)$ entries of the sketch \mathbf{S}_t at each iteration. Meanwhile, CGAL must update all $\Theta(n^2)$ entries of the matrix iterate \mathbf{X}_t .

As a consequence of these improvements, *SketchyCGAL* scales to problems that are orders of magnitude larger than CGAL can handle. See section 7.

7. Numerical examples. This section showcases computational experiments that establish that *SketchyCGAL* is a practical method for solving large SDPs. We show that the algorithm is **flexible** by applying it to several classes of SDPs, and we show it is **reliable** by solving a large number of instances of each type. We give empirical evidence that the (implicit) iterates are **converging** to optimality much faster than Theorem 6.3 suggests. Explicit comparisons with other general SDP solvers demonstrate that *SketchyCGAL* is **competitive** for small SDPs, while it **scales** to problems that standard methods cannot handle.

7.1. Setup. We begin with an overview of the experimental regimen.

7.1.1. Computational environment. All experiments are performed in MATLAB_R2018a with double-precision arithmetic (8 bytes per real float). Source code is included with the supplementary material. To simulate the processing power of a personal computer, we use a single Intel Xeon CPU E5-2630 v3, clocked at 2.40 GHz, with RAM usage capped at 16 GB.

Arithmetic costs are measured in terms of actual CPU time. MATLAB does not currently offer a memory profiler, so we externally monitor the total memory allocated. We approximate the storage cost by reporting the peak value minus the storage at MATLAB’s idle state.

7.1.2. Problem scaling. The bounds in the convergence theorem [Fact 3.1](#) for the implicit iterates of SketchyCGAL depend on problem scaling. Our analysis motivates us to set

$$(7.1) \quad \|C\|_F = \|\mathcal{A}\| = \alpha = 1 \quad \text{and} \quad \|A_1\|_F = \|A_2\|_F = \cdots = \|A_d\|_F.$$

In our experiments, we enforce the scalings [\(7.1\)](#), except where noted.

7.1.3. Quality of solutions. Given a prospective solution X of [\(2.2\)](#), we compute the relative suboptimality and feasibility as

$$\text{objective residual} = \frac{|\langle C, X \rangle - \langle C, X_\star \rangle|}{\max\{1, |\langle C, X_\star \rangle|\}} \quad \text{and} \quad \text{infeasibility} = \frac{\|\mathcal{A}X - b\|}{\max\{1, \|b\|\}}.$$

It is standard to place the maximum in each denominator to handle small values gracefully. These quantities are evaluated with respect to the original (not rescaled) problem data.

7.1.4. Implementation. Our experiments require a variant of SketchyCGAL that can handle inequality constraints; see [Appendix D](#). We implement the algorithm with the default parameters ($\beta_0 = 1$ and $K = +\infty$), and the sketch uses a Gaussian test matrix. The eigenvalue subproblem is solved via randomized Lanczos ([Algorithm 4.2](#)). We include the optional trace normalization (line 13 in [Algorithm 6.1](#)) whenever appropriate. No further tuning is done.

7.2. The MaxCut SDP. We begin with MaxCut SDPs [\(1.3\)](#). Our goal is to assess the storage and arithmetic costs of SketchyCGAL for a standard testbed. We compare with provable solvers for general SDPs: SEDUMI [\[97\]](#), SDPT3 [\[99\]](#), MOSEK [\[78\]](#), and SDPNAL+ [\[108\]](#).

7.2.1. Rounding. To extract a cut from an approximate solution to [\(1.3\)](#), we apply a simple rounding procedure. SketchyCGAL returns a matrix $\widehat{X} = U\Lambda U^*$, where $U \in \mathbb{R}^{n \times R}$ has orthonormal columns. The columns of the entrywise signum, $\text{sgn}(U)$, are the signed indicators of R cuts. We compute the weights of all R cuts and select the largest. The other solvers return a full-dimensional solution \widehat{X} ; we compute the top R eigenvectors of \widehat{X} using the MATLAB command `eigs` and invoke the same rounding procedure.

7.2.2. Datasets. We consider datasets from two different benchmark groups:

1. GSET: 67 binary-valued matrices generated by an autonomous random graph generator and published online [\[109\]](#). The dimension n varies from 800 to 10 000.
2. DIMACS10: This benchmark consists of 152 symmetric matrices (with n varying from 39 to 50 912 018) chosen for the 10th DIMACS Implementation Challenge [\[13\]](#). We consider the 138 datasets with dimension $n \leq 8\,000\,000$.

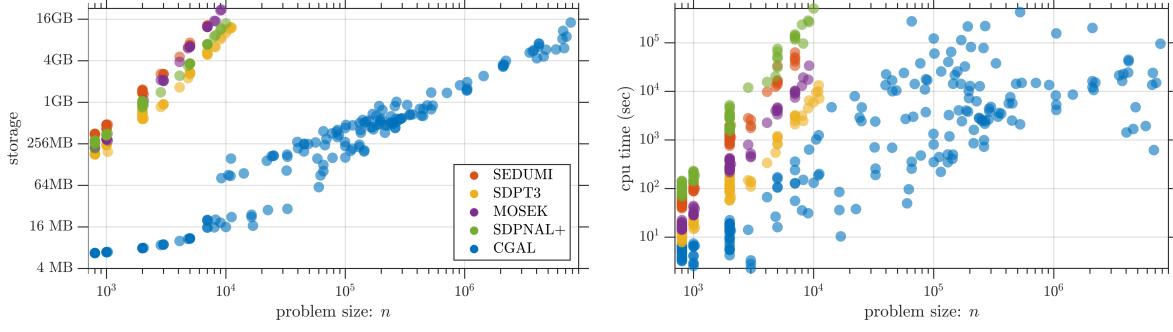


Figure 7.1. MaxCut SDP: Scalability. Storage cost [left] and runtime [right] of *SketchyCGAL* with sketch size $R = 10$ as compared with four standard SDP solvers. Each marker is one dataset. See [subsection 7.2.3](#).

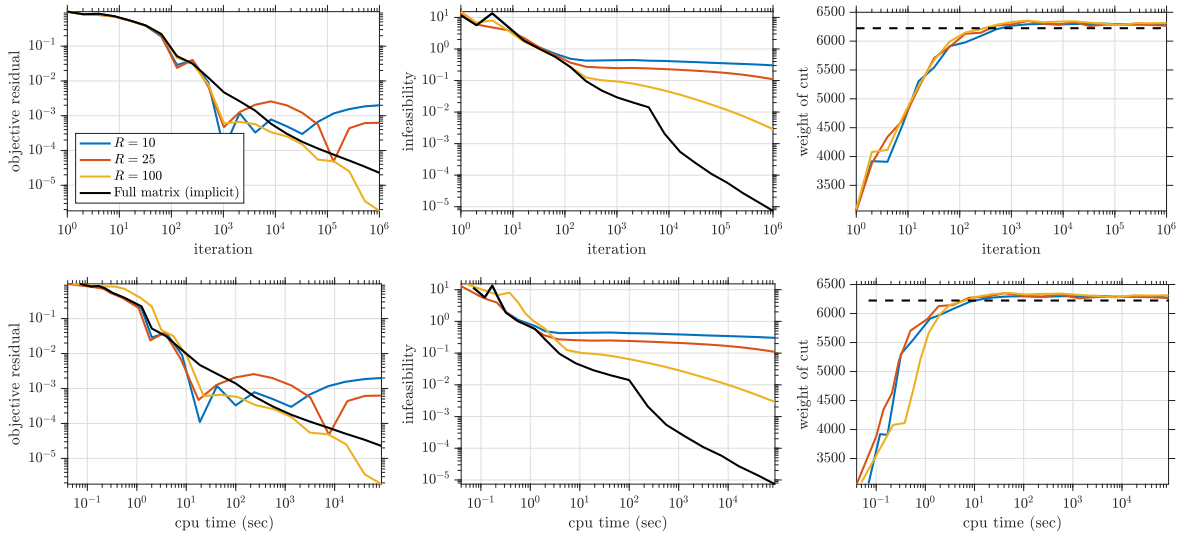


Figure 7.2. MaxCut SDP: Convergence. We solve the MaxCut SDP for the G67 dataset ($n = 10\,000$) with *SketchyCGAL*. The subplots show the suboptimality [left], infeasibility [center], and cut value [right] of the implicit iterates and low-rank approximations for sketch size $R \in \{10, 25, 100\}$ as a function of iteration [top] and CPU time [bottom]. The dashed line is the cut value of a high-accuracy SDPT3 solution. See [subsection 7.2.4](#).

7.2.3. Storage and arithmetic comparisons. We run each solver for each dataset and measure the storage cost and the runtime. We invoke *SketchyCGAL* with rank parameter $R = 10$, and we stop the algorithm when the error bound described in [Appendix C.1](#) guarantees that the implicit iterates have both suboptimality and feasibility below 10^{-4} . We use other solvers with their default parameter settings and stopping criteria.

[Figure 7.1](#) displays the results of this experiment. The conventional convex solvers do not scale beyond $n = 10^4$ due to their high storage costs. In contrast, *SketchyCGAL* handles problems with $n \approx 8 \cdot 10^6$. Further results appear in [Appendix F.1](#).

7.2.4. Empirical convergence rates. Next, we investigate the empirical convergence of *SketchyCGAL* and the effect of the sketch size parameter R . We consider the MaxCut SDP (1.3) for the G67 dataset ($n = 10\,000$), the largest instance in GSET, and we use a high-accuracy

solution from SDPT3 to approximate the optimal point. We run 10^6 iterations of SketchyCGAL for each $R \in \{10, 25, 100\}$.

Figure 7.2 illustrates the convergence of SketchyCGAL for this problem. After t iterations, the residual and infeasibility decline as $O(t^{-1})$, which is far better than the $O(t^{-1/2})$ bound in Theorem 6.3. Similar behavior is manifest for other datasets and other problems.

Figure 7.2 also displays the weight of the cut obtained after rounding, compared with the weight of the cut obtained from the SDPT3 solution. Observe that sketch size $R = 10$ is sufficient, and the SketchyCGAL solutions yield excellent cuts after a few hundred iterations.

7.2.5. Hard MaxCut instances. Waldspurger & Waters [105] construct instances of the MaxCut SDP (1.3) that are challenging for algorithms based on the Burer–Monteiro [27] factorization method (subsection 8.6). Each instance has a unique solution and the solution has rank 1, but Burer–Monteiro methods require factorization rank $R = \Theta(\sqrt{n})$, resulting in storage cost $\Theta(n^{3/2})$. In Appendix F.2, we give numerical evidence that SketchyCGAL solves these instances with sketch size $R = 2$, achieving the optimal storage $\Theta(n)$.

7.3. Abstract phase retrieval. Phase retrieval is the problem of reconstructing a complex-valued signal from intensity-only measurements. It arises in interferometry [41], speech processing [16], array imaging [33], microscopy [55], and many other applications. We will outline a standard method [33, 55] for performing phase retrieval by means of an SDP.

This section uses synthetic instances of a phase retrieval SDP to compare the scaling behavior of SketchyCGAL and CGAL. We also consider a third algorithm ThinCGAL, inspired by [112], that maintains a thin SVD of the matrix variable via rank-one updates [25].

7.3.1. Phase retrieval SDPs. Let $\mathbf{x}_{\mathfrak{h}} \in \mathbb{C}^n$ be an unknown (discrete) signal. For known vectors $\mathbf{a}_i \in \mathbb{C}^n$, we acquire measurements of the form

$$(7.2) \quad b_i = |\langle \mathbf{a}_i, \mathbf{x}_{\mathfrak{h}} \rangle|^2 \quad \text{for } i = 1, 2, 3, \dots, d.$$

Abstract phase retrieval is the challenging problem of recovering $\mathbf{x}_{\mathfrak{h}}$ from \mathbf{b} .

Let us summarize a lifting approach introduced by Balan et al. [15]. The key idea is to replace the signal vector $\mathbf{x}_{\mathfrak{h}}$ by the matrix $\mathbf{X}_{\mathfrak{h}} = \mathbf{x}_{\mathfrak{h}}\mathbf{x}_{\mathfrak{h}}^*$. Then rewrite (7.2) as

$$(7.3) \quad b_i = \mathbf{a}_i^* \mathbf{X}_{\mathfrak{h}} \mathbf{a}_i = \langle \mathbf{A}_i, \mathbf{X}_{\mathfrak{h}} \rangle \quad \text{where } \mathbf{A}_i = \mathbf{a}_i \mathbf{a}_i^* \quad \text{for } i = 1, 2, 3, \dots, d.$$

Promoting the implicit constraints on $\mathbf{X}_{\mathfrak{h}}$ and forming the \mathbf{A}_i into a linear map \mathcal{A} , we can express the problem of finding $\mathbf{X}_{\mathfrak{h}}$ as a feasibility problem with a matrix variable:

$$\text{find } \mathbf{X} \in \mathbb{S}_n \quad \text{subject to} \quad \mathcal{A}\mathbf{X} = \mathbf{b}, \quad \mathbf{X} \text{ is psd}, \quad \text{rank}(\mathbf{X}) = 1.$$

To reach a tractable convex formulation, we pass to a trace minimization SDP [40]:

$$(7.4) \quad \text{minimize} \quad \text{tr } \mathbf{X} \quad \text{subject to} \quad \mathcal{A}\mathbf{X} = \mathbf{b}, \quad \mathbf{X} \text{ is psd}, \quad \text{tr } \mathbf{X} \leq \alpha.$$

The parameter α is an upper bound on the signal energy $\|\mathbf{x}_{\mathfrak{h}}\|^2$, which can be estimated from the observed data \mathbf{b} ; see [112] for the details. We can solve the SDP (7.4) via SketchyCGAL.

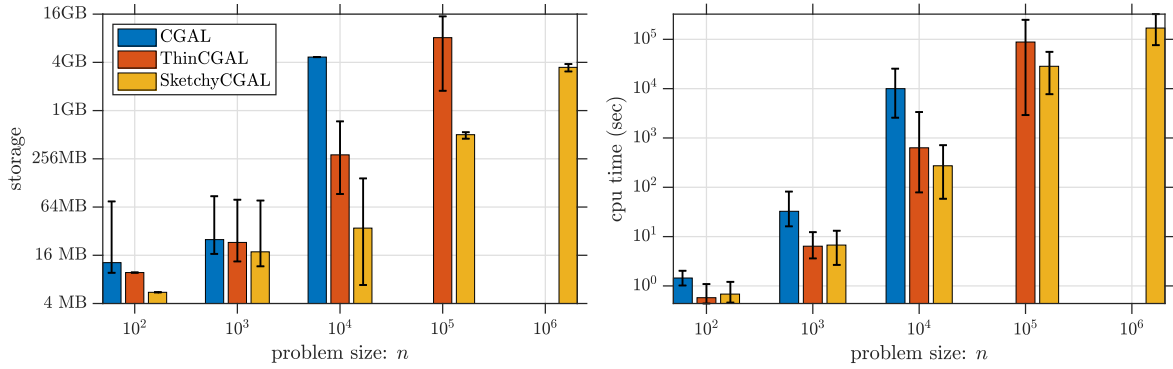


Figure 7.3. Phase retrieval SDP: Scalability. Storage cost [left] and runtime [right] to solve random instances with algorithms CGAL, ThinCGAL, and SketchyCGAL. Height of the bar is the mean and the interval marks are the min/max over 20 trials. Missing bars indicate the failure. See subsection 7.3.4.

7.3.2. Rounding. Suppose that we have obtained an approximate solution \mathbf{X} to (7.4). To estimate the signal $\chi_{\mathfrak{h}}$, we form the vector $\chi = \sqrt{\lambda} \mathbf{u}$ where (λ, \mathbf{u}) is a maximum eigenpair of \mathbf{X} . Both SketchyCGAL and ThinCGAL return eigenvalue decompositions, so this step is trivial. For CGAL, we use the MATLAB function `eigs` to perform this computation.

7.3.3. Dataset and evaluation. We consider synthetic phase retrieval instances. For each $n \in \{10^2, 10^3, \dots, 10^6\}$, we generate 20 independent datasets as follows. First, draw $\chi_{\mathfrak{h}} \in \mathbb{C}^n$ from the complex standard normal distribution. Then acquire $d = 10n$ phaseless measurements (7.2) using the coded diffraction pattern model [31]; see Appendix E.1. The induced linear maps \mathcal{A} and \mathcal{A}^* can be applied via the fast Fourier transform (FFT).

The relative error in a signal reconstruction χ is given by $\min_{\phi \in \mathbb{R}} \|e^{i\phi} \chi - \chi_{\mathfrak{h}}\| / \|\chi_{\mathfrak{h}}\|$. In the SDP (7.4), we set $\alpha = 3n$ to demonstrate insensitivity of the algorithm to the parameter.

7.3.4. Storage and arithmetic comparisons. For each algorithm, we report the storage cost and runtime required to produce a signal estimate χ with (exact) relative error below 10^{-2} . We invoke SketchyCGAL with sketch size parameter $R = 5$.

Figure 7.3 displays the outcome. We witness the benefit of sketching for both storage and arithmetic. CGAL fails for all large-scale instances ($n = 10^5$ and 10^6) due to storage allocation. For the same reason, ThinCGAL solves 19 of 20 instances for $n = 10^5$; it always fails for $n = 10^6$. SketchyCGAL successfully solves all problem instances to the target accuracy.

7.4. Phase retrieval in microscopy. Next, we study a more realistic phase retrieval problem that arises from a type of microscopy system [119] called Fourier ptychography (FP). Phase retrieval SDPs offer a potential approach to FP imaging [55]. This section shows that SketchyCGAL can successfully solve the difficult phase retrieval SDPs that arise from FP.

7.4.1. Fourier ptychography. FP microscopes circumvent the physical limits of a simple lens to achieve high-resolution and wide field-of-view simultaneously [119]. To do so, an FP microscope illuminates a sample from many angles and uses a simple lens to collect low-resolution intensity-only images. The measurements are low-pass filters, whose transfer functions depend on the lens and the angle of illumination [55]. From the data, we form a

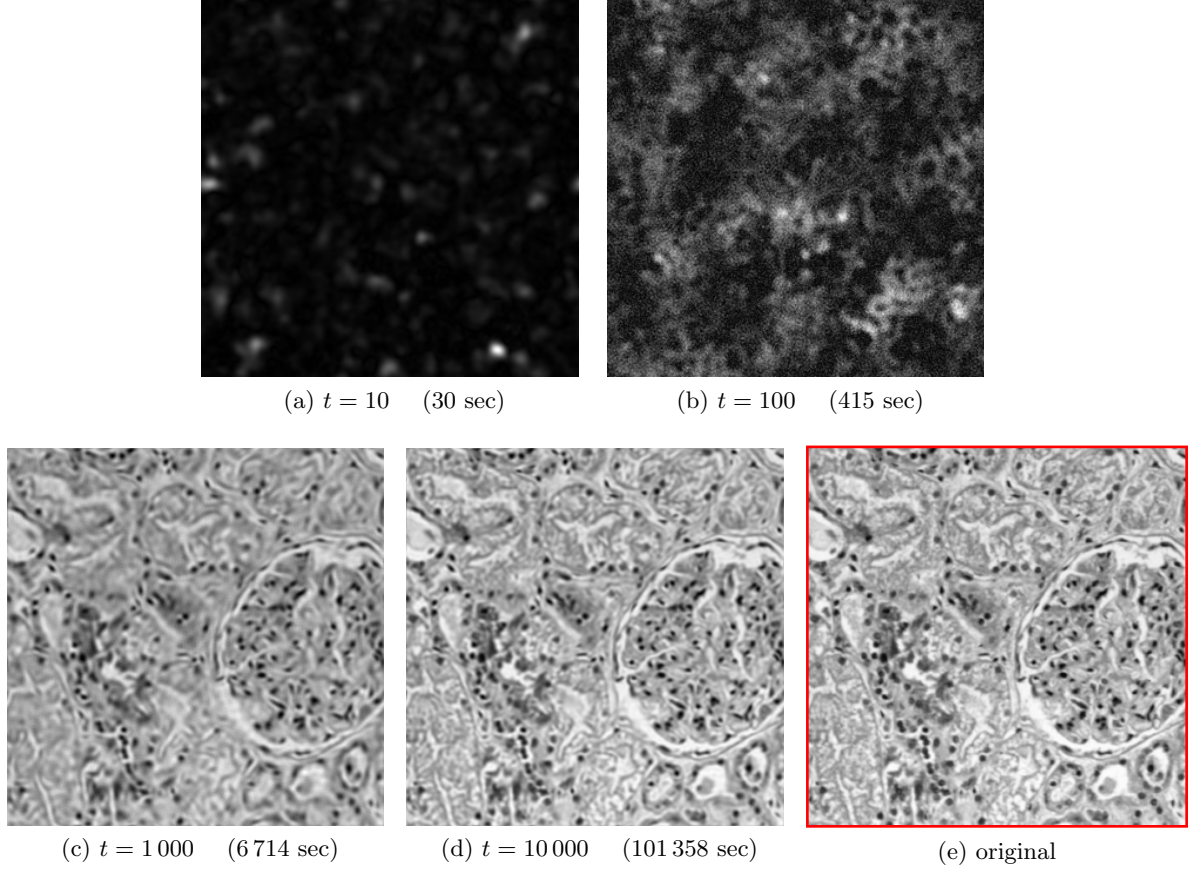


Figure 7.4. Phase retrieval SDP: Imaging. *Reconstruction of an $n = 320^2$ pixel image from Fourier ptychography data. We solve an $n \times n$ phase retrieval SDP via *SketchyCGAL* with rank parameter $R = 5$ and show the images obtained at iterations $t = 10, 10^2, 10^3, 10^4$. The last subfigure is the original. See [subsection 7.4.3](#).*

high-resolution image by solving a phase retrieval problem; e.g., via the SDP (7.4).

The high-resolution image of the sample is represented by a Fourier-domain vector $\chi_{\mathfrak{h}} \in \mathbb{C}^n$. We acquire d intensity-only measurements of the form (7.2), where d is the total number of pixels in the low-resolution illuminations. The low-pass measurements are encoded in vectors \mathbf{a}_i . The operators \mathcal{A} and \mathcal{A}^* , built from the matrices $\mathbf{A}_i = \mathbf{a}_i \mathbf{a}_i^*$, can be applied via the FFT.

7.4.2. Dataset and evaluation. The authors of [55] provided transmission matrices \mathbf{A}_i of a working FP system. We simulate this system in the computer environment to acquire noiseless intensity-only measurements of a high-resolution target image [36]. In this setup, $\chi_{\mathfrak{h}} \in \mathbb{C}^n$ corresponds to the Fourier transform of an $n = 320^2 = 102\,400$ pixel grayscale image. We normalize $\chi_{\mathfrak{h}}$ so that $\|\chi_{\mathfrak{h}}\| = 1$. We acquire 225 low-resolution illuminations of the original image, each with $64^2 = 4\,096$ pixels. The total number of measurements is $d = 921\,600$.

We evaluate the error of an estimate χ as in [subsection 7.3.3](#). Although we know that the signal energy equals 1, it is more realistic to approximate the signal energy by setting $\alpha = 1.5$ in the SDP (7.4).

7.4.3. FP imaging. We solve the phase retrieval SDP (7.4) by performing 10 000 iterations of SketchyCGAL with rank parameter $R = 5$. The top eigenvector of the output gives an approximation $\chi \in \mathbb{C}^n$ of the signal. The inverse Fourier transform of χ is the desired image.

Figure 7.4 displays the image reconstruction after $t \in \{10, 10^2, 10^3, 10^4\}$ iterations. We obtain a good-quality result in 2 hours after 1 000 iterations; a sharper image emerges in 28 hours after 10 000 iterations are complete. We believe the computational time can be reduced substantially with a parallel or GPU implementation. Nevertheless, it is gratifying that we have solved an SDP whose matrix variable has over 10^{10} entries! See Appendix E.2 for a larger FP instance with $n = 640^2$ pixels.

7.5. The quadratic assignment problem. The quadratic assignment problem (QAP) is a very difficult combinatorial optimization problem that includes the traveling salesman, max-clique, bandwidth problems, and many others as special cases [74]. SDP relaxations offer a powerful approach for obtaining good solutions to large QAP problems [117]. In this section, we demonstrate that SketchyCGAL can solve these challenging SDPs.

7.5.1. QAP. We begin with the simplest form of the QAP. Fix symmetric $\mathbf{n} \times \mathbf{n}$ matrices $\mathbf{A}, \mathbf{B} \in \mathbb{S}_{\mathbf{n}}$ where \mathbf{n} is a natural number. We wish to “align” the matrices by solving

$$(7.5) \quad \text{minimize} \quad \text{tr}(\mathbf{A}\mathbf{\Pi}\mathbf{B}\mathbf{\Pi}^*) \quad \text{subject to} \quad \mathbf{\Pi} \text{ is an } \mathbf{n} \times \mathbf{n} \text{ permutation matrix.}$$

Recall that a permutation matrix $\mathbf{\Pi}$ has precisely one nonzero entry in each row and column, and that nonzero entry equals one. A brute force search over the $\mathbf{n}!$ permutation matrices of size \mathbf{n} quickly becomes intractable as \mathbf{n} grows. Unsurprisingly, the QAP problem (7.5) is NP-hard [94]. Instances with $\mathbf{n} > 30$ usually cannot be solved in reasonable time.

7.5.2. Relaxations. There is an extensive literature on SDP relaxations for QAP, beginning with the work [117] of Zhao et al. We consider a weaker relaxation inspired by [56, 26]:

$$(7.6) \quad \begin{aligned} & \text{minimize} \quad \text{tr}[(\mathbf{B} \otimes \mathbf{A})\mathbf{Y}] \\ & \text{subject to} \quad \text{tr}_1(\mathbf{Y}) = \mathbf{I}, \quad \text{tr}_2(\mathbf{Y}) = \mathbf{I}, \quad \mathcal{G}(\mathbf{Y}) \geq 0, \\ & \quad \text{vec}(\mathbf{P}) = \text{diag}(\mathbf{Y}), \quad \mathbf{P}\mathbf{1} = \mathbf{1}, \quad \mathbf{1}^*\mathbf{P} = \mathbf{1}^*, \quad \mathbf{P} \geq \mathbf{0}, \\ & \quad \begin{bmatrix} 1 & \text{vec}(\mathbf{P})^* \\ \text{vec}(\mathbf{P}) & \mathbf{Y} \end{bmatrix} \succcurlyeq \mathbf{0}, \quad \text{tr } \mathbf{Y} = n. \end{aligned}$$

We have written \otimes for the Kronecker product.

The constraint $\mathcal{G}(\mathbf{Y}) \geq 0$ enforces nonnegativity of a subset of the entries in \mathbf{Y} . In the Zhao et al. relaxation, \mathcal{G} is the identity map, so it yields $O(\mathbf{n}^4)$ constraints. We reduce the complexity by choosing \mathcal{G} more carefully. In our formulation, \mathcal{G} extracts precisely the nonzero entries of $\mathbf{B} \otimes \mathbf{1}\mathbf{1}^*$. This is beneficial because \mathbf{B} is sparse in many applications. For example, in traveling salesman and bandwidth problems, \mathbf{B} has $O(\mathbf{n})$ nonzero entries, so the map \mathcal{G} produces only $O(\mathbf{n}^3)$ constraints.

The main variable \mathbf{Y} in (7.6) has dimension $\mathbf{n}^2 \times \mathbf{n}^2$. As a consequence, the problem has $O(\mathbf{n}^4)$ degrees of freedom, together with $O(\mathbf{n}^3)$ to $O(\mathbf{n}^4)$ constraints (depending on \mathcal{G}). The explosive growth of this relaxation scuttles most algorithms by the time $\mathbf{n} > 50$. To solve larger instances, many researchers resort to even weaker relaxations.

In contrast, we can solve the relaxation (7.6) directly using SketchyCGAL, up to $n = 150$. By limiting the number of inequality constraints, via the operator \mathcal{G} , we achieve substantial reductions in resource usage. We validate our algorithm on QAPs where the exact solution is known, and we compare the performance with algorithms [116, 26] for other relaxations.

7.5.3. Rounding. Given an approximate solution to (7.6), we use a rounding method to construct a permutation. SketchyCGAL returns a matrix $\widehat{\mathbf{X}} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$ where \mathbf{U} has dimension $(n^2 + 1) \times R$. We extract the first column of \mathbf{U} , discard its first entry and reshape the remaining part into an $n \times n$ matrix. Then we project this matrix onto the set of permutation matrices via the Hungarian method [66, 79, 59]. This yields a feasible point $\mathbf{\Pi}$ for the problem (7.5). We repeat this procedure for all R columns of \mathbf{U} , and we pick the one that minimizes $\text{tr}(\mathbf{A}\mathbf{\Pi}\mathbf{B}\mathbf{\Pi}^*)$. This permutation gives an upper bound on the optimal value of (7.5).

7.5.4. Datasets and evaluation. We consider instances from QAPLIB [30] and TSPLIB [91] that are used in [26]. The optimal values are known, and the permutation size n varies between 12 and 150. (Recall that the SDP matrix dimension $n = n^2 + 1$.) We report

$$(7.7) \quad \text{relative gap \%} = \frac{\text{upper bound obtained} - \text{optimum}}{\text{optimum}} \times 100$$

7.5.5. Solving QAPs. To solve (7.6), we cannot use the scaling (7.1) because $\|\mathcal{A}\|$ is not available; see the source code for our approach. We apply SketchyCGAL with sketch size $R = n$, using total storage $\Theta(n^3)$. After rounding, a low- or medium-accuracy solution of (7.6) often provides a better permutation than a high-accuracy solution. Therefore, we applied the rounding step at iterations 2, 4, 8, 16, \dots and tracked the quality of the best permutation attained on the solution path. We stopped after (the first of) 10^6 iterations or 48 hours.

The results of this experiment appear in Figure 7.5. We compare against the best value reported by Bravo Ferreira et al. in [26, Tables 4 and 6] for their CSDP method with clique size $k \in \{2, 3, 4\}$. We also include the results that [26] lists for the PATH method [116].

SketchyCGAL allows us to solve a tighter SDP relaxation of QAP than the other methods (CSDP, PATH). As a consequence, we obtain significantly smaller gaps for most instances. In fact, SketchyCGAL even recovers the exact solution for some of the instances in QAPLIB!

8. Related work. There is a large body of literature on numerical methods for solving SDPs. Most of these algorithms adopt standard techniques from optimization. This section describes the major approaches, but an exhaustive discussion is beyond the scope of this work. See [75] for a recent survey.

8.1. Interior point method (IPM). The first practical IPM algorithms for linear programming were developed in 1984 by Karmarkar [62]. Soon after, they were extended to SDPs independently by Nesterov & Nemirovski [83, 84], Alizadeh [3, 4], and Karmarkar [60, 61]. IPMs are still widely used for solving SDPs to high precision. Software packages include SeDuMi [97], MoSeK [78], and SDPT3 [99].

Unfortunately, IPMs do not scale to large problems. Indeed, an IPM formulates a sequence of unconstrained barrier problems, which are solved via Newton’s method or its variants. This step requires storing and factoring large, dense matrices. As a result, a typical IPM for the SDP (2.2) requires $\Theta(n^3 + d^2n^2 + d^3)$ arithmetic operations per iteration and $\Theta(n^2 + dn + d^2)$

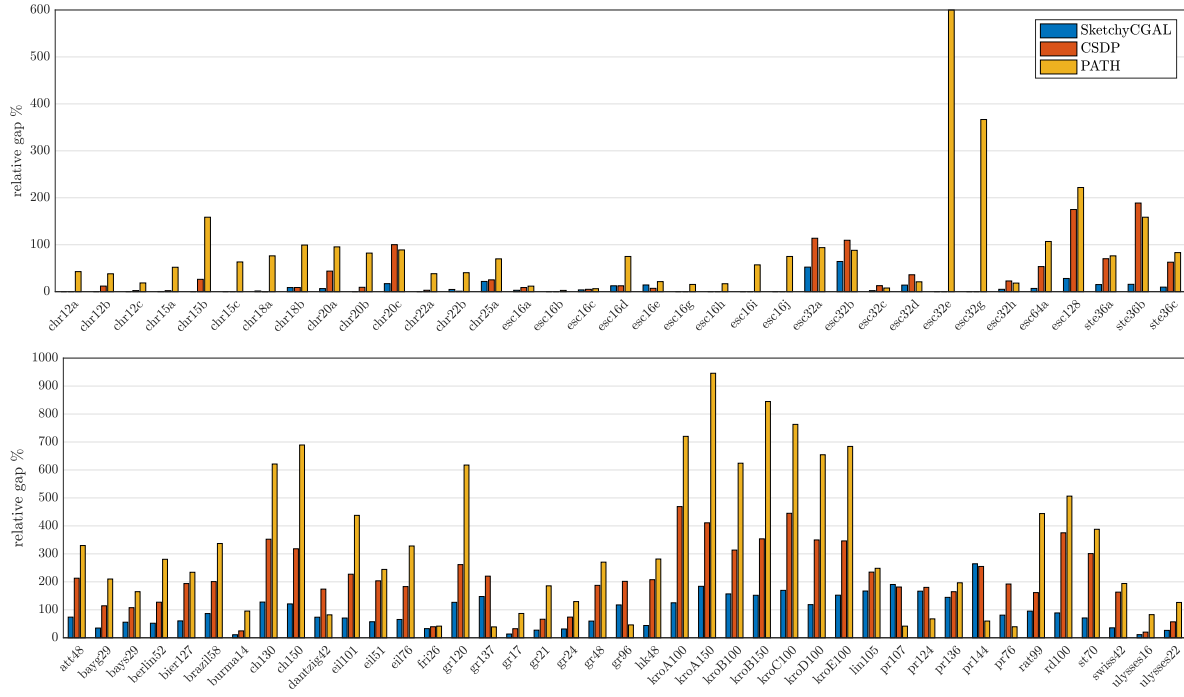


Figure 7.5. QAP relaxation: Solution quality. Using *SketchyCGAL*, the *CSDP* method [26], and the *PATH* method [116], we solve SDP relaxations of QAP instances from QAPLIB [top] and TSPLIB [bottom]. The bars compare the cost of the computed solution, against the optimal value; shorter is better. See subsection 7.5.5.

memory [6]. As d and n grow, time and storage costs become prohibitive, even for moderately sized problems, such as $n, d \gtrsim 10\,000$. For more details, see [64, 9, 81, 19].

8.2. Augmented Lagrangian methods. Several effective SDP solvers are based on the augmented Lagrangian (AL) paradigm. These methods differ in the ways that the constraints are formulated and the subproblems are solved. See [29, 106, 118, 107] for examples.

Although AL techniques are often called “first-order methods,” some of these algorithms use second-order information. In particular, Zhao et al. [118] employ a semi-smooth Newton method and the conjugate gradient method to solve the subproblems. Their method is enhanced and implemented in the software package SDPNAL+ [108].

Classical AL methods for SDPs typically require storage $\Omega(n^2)$. To address this weakness, researchers have combined the AL framework with CGM [45, 73, 95, 110] and with nonconvex matrix factorization [27, 28, 67, 93]. We explain some of these ideas below.

8.3. First-order methods. Next, we summarize the first-order methods for large SDPs. We focus on *projection-free* algorithms, which do not require full SVD computations.

8.3.1. Bundle methods. One of the earliest first-order SDP algorithms is the spectral bundle method developed by Helmberg and coauthors [53, 52, 51]. These algorithms build a piecewise affine model for the optimization problem using eigenvectors of the dual slack matrix $C - \mathcal{A}y^*$. At present, these methods lack a complete theoretical analysis.

8.3.2. Conditional gradient method (CGM) and friends. CGM solves a smooth convex minimization problem by repeatedly optimizing a linear approximation for the objective function over the constraint set. This approach was introduced by Frank & Wolfe [42] for optimization on polytopes; see also [71, 58]. Much later, Clarkson [35] developed a new analysis for optimization on the probability simplex Δ_n , and Hazan [50] studied optimization on the spectahedron Δ_n . Jaggi [57] generalized CGM to all compact, convex constraint sets.

The standard CGM algorithm does not apply to the model problem (2.2) because of the affine constraint $\mathcal{A}\mathbf{X} = \mathbf{b}$. Several variants [45, 73, 95, 111] of CGM can handle affine constraints. In particular, CGAL [110] does so by applying CGM to an AL formulation. We have chosen to extend CGAL because of its strong empirical performance; see [110, Sec. 4-5].

8.3.3. Primal-dual subgradient methods. Several authors have developed primal-dual subgradient algorithms [82, 113] that can be applied to SDPs. These methods perform subgradient ascent on the dual problem; the cost of each iteration is dominated by an eigenvector computation. Nesterov [82] constructs primal iterates by proximal mapping. The algorithm in Yurtsever et al. [113] constructs primal iterates by an averaging technique, similar to the updates in CGM. Bach [12] has described a relationship between CGM and dual subgradient ascent. In fact, CGAL and related methods are types of primal-dual averaging [114, 110].

8.4. Primal-dual games and online learning methods. Motivated by applications in theoretical computer science, researchers in this field have also devised SDP algorithms.

This approach first reduces SDP optimization to a sequence of feasibility problems. By standard duality techniques, each feasibility problem admits a saddle-point formulation:

$$(8.1) \quad \max_{\mathbf{X} \in \alpha \Delta_n} \min_{\mathbf{p} \in \Delta_d} \langle \mathbf{p}, \mathcal{A}\mathbf{X} - \mathbf{b} \rangle,$$

where Δ_d denotes the $(d-1)$ -dimensional probability simplex. The dual of (8.1) is an eigenvalue optimization problem, which is important in its own right.

Most theoretical approaches to (8.1) are based on the matrix multiplicative weight (MMW) method [103, 10], which is nothing but the mirror-prox algorithm with the entropy mirror map [80]. The basic idea is to perform gradient descent in a dual space, using the matrix exponential map to transfer information back to the primal space. To make this approach more scalable, researchers have proposed linearization, random projection, and sparsification techniques to approximate the matrix exponential; see [11, 7, 88, 43, 44, 8, 14, 69].

This line of work has culminated in a simple rank-one MMW method [32], which uses the stochastic Lanczos method to approximate the matrix exponential. Even so, the reduction to a sequence of feasibility problems makes this technique impractical for solving general SDPs.

We are aware of only one experimental evaluation of the MMW idea, which appears in the randomized mirror-prox paper of Baes et al. [14].

8.5. Storage considerations. Almost all provably correct SDP algorithms store and operate on a full-dimensional matrix variable, so they are not suitable for large-scale SDPs.

Some primal-dual subgradient methods and CGM variants build an approximate solution as a convex combination of rank-one updates, so the rank of the solution does not exceed the number of iterations. This fact has led researchers to call these methods “storage-efficient,” but this claim is misleading because the algorithms require many iterations to converge.

Recently, a subset of the authors has shown how to combine sketching with CGM to design a storage-optimal algorithm, **SketchyCGM**, for a class of low-rank matrix optimization problems [115]. We believe that **SketchyCGM** was the first provably correct storage-optimal method for these problems. **SketchyCGAL** was inspired by **SketchyCGM**.

In concurrent work with Ding [39], we developed another storage-optimal algorithm for standard-form SDPs. This method is based on a new *approximate complementarity principle*. Roughly, we can use a suboptimal dual point to approximate the range of the primal solution to the SDP. By compressing the primal problem to this subspace, we can solve the SDP with limited storage. This method, however, may fail on a set of SDPs with measure zero, and it has more limited guarantees than **SketchyCGAL**. A numerical evaluation is in process.

We have also noticed that the rank-one MMW method [32] can be combined with the Nyström sketch (Algorithm 5.1) to obtain a simple, storage-optimal algorithm for (8.1).

8.6. Nonconvex Burer–Monteiro methods. The most famous approach to low-storage semidefinite programming is the factorization heuristic proposed by Homer & Peinado [54] and extended by Burer & Monteiro (BM) [27]. The main idea is to rewrite the psd matrix variable $\mathbf{X} = \mathbf{F}\mathbf{F}^*$ in terms of a factor $\mathbf{F} \in \mathbb{R}^{n \times R}$, where the rank parameter $R \ll n$. We can reformulate the model problem (2.2) in terms of the new variable:

$$(8.2) \quad \text{minimize } \langle \mathbf{C}, \mathbf{F}\mathbf{F}^* \rangle \quad \text{subject to } \mathcal{A}(\mathbf{F}\mathbf{F}^*) = \mathbf{b}, \quad \mathbf{F}\mathbf{F}^* \in \alpha \Delta_n.$$

This approach controls storage by sacrificing convexity and the associated guarantees. We can apply various nonlinear programming methods to optimize (8.2). AL methods are most commonly used [27, 28, 67, 93] because of their empirical success.

There has been an intense effort to establish theoretical guarantees for the BM factorization approach. It is clear that every solution to the SDP (2.2) of rank R or less is also a solution to the factorized problem (8.2). On the other hand, (8.2) may admit spurious local minima. Furthermore, because the formulation is nonconvex, we can only expect a numerical algorithm to locate a first- or second-order critical point.

Burer & Monteiro [28] developed the initial theory about (8.2): If \mathbf{F} is a local minimum of (8.2) and if $R \geq \sqrt{2(d+1)}$, then \mathbf{X} is either a global minimum or it is contained in the relative interior of a face of the feasible set of (2.2). Boumal et al. [23] obtained an improvement: If the constraint set of (8.2) is a smooth manifold and $R \geq \sqrt{2(d+1)}$ and \mathbf{C} is generic, then each second-order critical point of (8.2) is a global optimum. A second-order critical point can be located using a Riemannian trust region method. See [21, 89, 34] for more extensions, but this theory requires opaque technical conditions that may not hold in practice.

The storage and arithmetic costs of solving (8.2) depend on the factorization rank R , as well as the optimization algorithm. Existing theoretical guarantees demand that $R = \Omega(\sqrt{d})$. In contrast, there are many applications where the SDP has a unique solution with *constant* rank. Thus, it is natural to ask whether we can reduce the rank parameter in the factorization. Waldspurger & Waters [105] resolved this question in the negative. Roughly, the BM formulation (8.2) can have spurious solutions unless $R = \Omega(\sqrt{d})$, and the bad problem instances form a set of positive measure. As a consequence, the Burer–Monteiro approach cannot support provably correct algorithms with storage costs better than $\Omega(n\sqrt{d})$. See Appendix F.2 for numerical evidence.

The most popular research software based on BM factorization is **Manopt** [22]. Solvers in **Manopt** implement manifold optimization algorithms, including Riemannian gradient and Riemannian trust region methods [1]. Consequently, **Manopt** is limited to problems where the factorized formulation (8.2) defines a smooth manifold.

8.7. Conclusion. We have presented a practical, new approach for solving standard-form SDPs at scale. Our algorithm combines a primal-dual optimization method with randomized linear algebra techniques to achieve unprecedented guarantees. We hope that our ideas lead to further algorithmic advances and support new applications of semidefinite programming.

Appendix A. Analysis of the CGAL Algorithm.

This section contains a complete analysis of the convergence of the CGAL algorithm and its arithmetic costs. For simplicity, we have specialized this presentation to the model problem that is the focus of this paper; many extensions are possible. The convergence results here are adapted from the initial paper [110] on the CGAL algorithm. The analysis of the arithmetic cost for the model problem is new.

A.1. The model problem. We focus on solving the optimization template

$$(A.1) \quad \text{minimize } \langle C, X \rangle \quad \text{subject to } \mathcal{A}X = b, \quad X \in \alpha\Delta_n.$$

The constraint set $\alpha\Delta_n$ consists of $n \times n$ psd matrices with trace α . The objective function depends on a matrix $C \in \mathbb{S}_n$. The linear constraints are determined by the linear map $\mathcal{A} : \mathbb{S}_n \rightarrow \mathbb{R}^d$ and the right-hand side vector $b \in \mathbb{R}^d$.

A.2. Elements of Lagrangian duality. Introduce the Lagrangian function

$$(A.2) \quad L(X; y) := \langle C, X \rangle + \langle y, \mathcal{A}X - b \rangle \quad \text{for } X \in \alpha\Delta_n \text{ and } y \in \mathbb{R}^d.$$

We assume that the Lagrangian admits at least one saddle point $(X_\star, y_\star) \in \alpha\Delta_n \times \mathbb{R}^d$:

$$(A.3) \quad L(X_\star, y) \leq L(X_\star, y_\star) \leq L(X; y_\star) \quad \text{for all } X \in \alpha\Delta_n \text{ and } y \in \mathbb{R}^d.$$

This hypothesis is guaranteed under Slater's condition. Write p_\star for the shared extremal value of the dual and primal problems:

$$(A.4) \quad \max_{y \in \mathbb{R}^d} \min_{X \in \alpha\Delta_n} L(X, y) = p_\star = \min_{X \in \alpha\Delta_n} \sup_{y \in \mathbb{R}^d} L(X, y).$$

In particular, note that $p_\star = \langle C, X_\star \rangle$.

A.3. The CGAL iteration. Let $\beta_0 > 0$ be an initial smoothing parameter. Fix a schedule for the step size parameter η_t and the smoothing parameter β_t :

$$(A.5) \quad \eta_t := \frac{2}{t+1} \quad \text{and} \quad \beta_t := \beta_0 \sqrt{t+1} \quad \text{for } t = 1, 2, 3, \dots$$

Define the augmented Lagrangian L_t with smoothing parameter β_t

$$(A.6) \quad L_t(X; y) := \langle C, X \rangle + \langle y, \mathcal{A}X - b \rangle + \frac{1}{2}\beta_t \|\mathcal{A}X - b\|^2.$$

The CGAL algorithm solves the model problem (A.1) by alternating between primal and dual update steps on (A.6), while increasing the smoothing parameter.

Fix the initial iterates

$$(A.7) \quad \mathbf{X}_1 = \mathbf{0} \in \mathbb{S}_n \quad \text{and} \quad \mathbf{y}_1 = \mathbf{0} \in \mathbb{R}^d.$$

At each iteration $t = 1, 2, 3, \dots$, we implicitly compute the partial derivative

$$(A.8) \quad \mathbf{D}_t := \partial_{\mathbf{X}} L_t(\mathbf{X}_t; \mathbf{y}_t) = \mathbf{C} + \mathcal{A}^*(\mathbf{y}_t + \beta_t(\mathcal{A}\mathbf{X}_t - \mathbf{b})).$$

Then we identify a primal update direction $\mathbf{H}_t \in \alpha\Delta_n$ that satisfies

$$(A.9) \quad \langle \mathbf{D}_t, \mathbf{H}_t \rangle \leq \min_{\mathbf{H} \in \alpha\Delta_n} \langle \mathbf{D}_t, \mathbf{H} \rangle + \frac{\alpha\beta_0}{\beta_t} \|\mathbf{D}_t\|.$$

In other words, the smoothing parameter β_t also controls the amount of inexactness we are willing to tolerate in the linear minimization at iteration t . We construct the next primal iterate via the rule

$$(A.10) \quad \mathbf{X}_{t+1} = \mathbf{X}_t + \eta_t(\mathbf{H}_t - \mathbf{X}_t) \in \alpha\Delta_n.$$

The dual update takes the form

$$(A.11) \quad \mathbf{y}_{t+1} = \mathbf{y}_t + \gamma_t(\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}).$$

We select the largest dual step size parameter $0 \leq \gamma_t \leq \beta_0$ that satisfies the condition

$$(A.12) \quad \gamma_t \|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\|^2 \leq \beta_t \eta_t^2 \alpha^2 \|\mathcal{A}\|^2.$$

The latter inequality always holds when $\gamma_t = 0$. We will also choose the dual step size to maintain a bounded travel condition:

$$(A.13) \quad \|\mathbf{y}_t\| \leq K.$$

If the bounded travel condition holds at iteration $t - 1$, then the choice $\gamma_t = 0$ ensures that it holds at iteration t . In practice, it is not necessary to enforce (A.13). The iteration continues until it reaches a maximum iteration count T_{\max} .

A.4. Solution quality. We can evaluate the quality of a given primal iterate \mathbf{X}_t with the information we have at hand. For simplicity, let us assume that we have solved the linear minimization (A.9) exactly at iteration t . That is, there is no error depending on $\|\mathbf{D}_t\|$. We can make this calculation periodically (say, for $t = 2^i$ for $i \in \mathbb{N}$) to decide whether to stop.

First, note that we have instant access to the infeasibility: $\mathcal{A}\mathbf{X}_t - \mathbf{b}$. Second, we can bound the suboptimality of the primal objective value. Introduce the surrogate duality gap

$$g_t(\mathbf{X}) := \max_{\mathbf{H} \in \alpha\Delta_n} \langle \mathbf{D}_t, \mathbf{X} - \mathbf{H} \rangle.$$

The quantity $g_t(\mathbf{X}_t)$ can be determined from the solution to (A.9). Then

$$\langle \mathbf{C}, \mathbf{X}_t \rangle - p_* \leq g_t(\mathbf{X}_t) - \langle \mathbf{y}_t, \mathcal{A}\mathbf{X}_t - \mathbf{b} \rangle.$$

See Appendix A.7 for discussion and a more refined estimate for the suboptimality.

Algorithm A.1 CGAL for the model problem (A.1)**Input:** Problem data for (A.1) implemented via the primitives (2.4); number T of iterations**Output:** Approximate solution \mathbf{X}_T to (2.2)**Recommendation:** Set $T \approx \varepsilon^{-1}$ to achieve ε -optimal solution (A.17)

```

1 function CGAL( $T$ )
2   Scale problem data (subsection 7.1.2)                                ▷ [opt] Recommended!
3    $\beta_0 \leftarrow 1$  and  $K \leftarrow +\infty$                                 ▷ Default parameters
4    $\mathbf{X} \leftarrow \mathbf{0}_{n \times n}$  and  $\mathbf{y} \leftarrow \mathbf{0}_d$ 
5   for  $t \leftarrow 1, 2, 3, \dots, T$  do
6      $\beta \leftarrow \beta_0 \sqrt{t+1}$  and  $\eta \leftarrow 2/(t+1)$ 
7      $(\xi, \mathbf{v}) \leftarrow \text{ApproxMinEvec}(\mathbf{C} + \mathcal{A}^*(\mathbf{y} + \beta(\mathcal{A}\mathbf{X} - \mathbf{b})); q_t)$ 
                                                    ▷ Use primitives (2.4) ❶❷!
                                                    ▷ With Algorithm 4.1, set  $q_t = 8t^{1/2} \log n$ 
                                                    ▷ With Algorithm 4.2, set  $q_t = t^{1/4} \log n$ 
8      $\mathbf{X} \leftarrow (1 - \eta) \mathbf{X} + \eta(\alpha \mathbf{v} \mathbf{v}^*)$ 
9      $\mathbf{y} \leftarrow \mathbf{y} + \gamma(\mathcal{A}\mathbf{X} - \mathbf{b})$                                 ▷ Step size  $\gamma$  satisfies (A.12) and (A.13)

```

A.5. Theoretical analysis of CGAL. We develop two results on the behavior of CGAL. The first concerns its convergence to optimality, and the second concerns the computational resource usage.

A.5.1. Convergence. The first result demonstrates that the CGAL algorithm always converges to a primal optimal solution of the model problem (A.1). This result is adapted from [110]; a complete proof appears below in Appendix A.6.

Theorem A.1 (CGAL: Convergence). *Define*

$$E := 6\alpha^2 \|\mathcal{A}\|^2 + \alpha(\|\mathbf{C}\| + K\|\mathcal{A}\|).$$

The primal iterates $\{\mathbf{X}_t : t = 2, 3, 4, \dots\}$ generated by the CGAL iteration satisfy the a priori bounds

$$(A.14) \quad \langle \mathbf{C}, \mathbf{X}_t \rangle - p_\star \leq \frac{2\beta_0 E}{\sqrt{t}} + K \cdot \|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|;$$

$$(A.15) \quad -\|\mathbf{y}_\star\| \cdot \|\mathcal{A}\mathbf{X}_t - \mathbf{b}\| \leq \langle \mathbf{C}, \mathbf{X}_t \rangle - p_\star;$$

$$(A.16) \quad \|\mathcal{A}\mathbf{X}_t - \mathbf{b}\| \leq \frac{2\beta_0^{-1}(K + \|\mathbf{y}_\star\|) + 2\sqrt{E}}{\sqrt{t}}.$$

The a priori bounds ensure that

$$(A.17) \quad \|\mathcal{A}\mathbf{X}_t - \mathbf{b}\| \leq \varepsilon \quad \text{and} \quad |\langle \mathbf{C}, \mathbf{X}_t \rangle - p_\star| \leq \varepsilon$$

within $O(\varepsilon^{-2})$ iterations. The big-O suppresses constants that depend on the problem data $(\alpha, \|\mathcal{A}\|, \|\mathbf{C}\|)$ and the algorithm parameters β_0 and K .

A.5.2. Problem scaling. Theorem A.1 indicates that it is valuable to scale the model problem (A.1) so that $\alpha = \|C\| = \|\mathcal{A}\| = 1$. In this case, a good choice for the smoothing parameter is $\beta_0 = 1$. Nevertheless, the algorithm converges, regardless of the scaling and regardless of the parameter choices β_0 and K . We use a slightly different scaling in practice; see subsection 7.1.2.

A.6. Proof of Theorem A.1. In this section, we establish the convergence guarantee stated in Theorem A.1.

A.6.1. Analysis of the primal update. The first steps in the proof address the role of the primal update rule (A.10). The arguments parallels the standard convergence analysis [57] of CGM, applied to the function

$$(A.18) \quad f_t(\mathbf{X}) := L_t(\mathbf{X}; \mathbf{y}_t) = \langle \mathbf{C}, \mathbf{X} \rangle + \langle \mathbf{y}_t, \mathcal{A}\mathbf{X} - \mathbf{b} \rangle + \frac{1}{2}\beta_t\|\mathcal{A}\mathbf{X} - \mathbf{b}\|^2.$$

Observe that the gradient $\nabla f_t(\mathbf{X}_t)$ coincides with the partial derivative (A.8).

To begin, we exploit the smoothness of f_t to control the change in its value at adjacent primal iterates. The function f_t is convex on \mathbb{S}_n , and its gradient has Lipschitz constant $\beta_t\|\mathcal{A}\|^2$, so

$$f_t(\mathbf{X}_+) - f_t(\mathbf{X}) \leq \langle \nabla f_t(\mathbf{X}), \mathbf{X}_+ - \mathbf{X} \rangle + \frac{1}{2}\beta_t\|\mathcal{A}\|^2\|\mathbf{X}_+ - \mathbf{X}\|_{\mathbb{F}}^2 \quad \text{for } \mathbf{X}, \mathbf{X}_+ \in \mathbb{S}_n.$$

In particular, with $\mathbf{X}_+ = \mathbf{X}_{t+1}$ and $\mathbf{X} = \mathbf{X}_t$, we obtain

$$(A.19) \quad \begin{aligned} f_t(\mathbf{X}_{t+1}) &\leq f_t(\mathbf{X}_t) + \langle \mathbf{D}_t, \mathbf{X}_{t+1} - \mathbf{X}_t \rangle + \frac{1}{2}\beta_t\|\mathcal{A}\|^2\|\mathbf{X}_{t+1} - \mathbf{X}_t\|_{\mathbb{F}}^2 \\ &= f_t(\mathbf{X}_t) + \eta_t \langle \mathbf{D}_t, \mathbf{H}_t - \mathbf{X}_t \rangle + \frac{1}{2}\beta_t\eta_t^2\|\mathcal{A}\|^2\|\mathbf{H}_t - \mathbf{X}_t\|_{\mathbb{F}}^2 \\ &\leq f_t(\mathbf{X}_t) + \eta_t \langle \mathbf{D}_t, \mathbf{H}_t - \mathbf{X}_t \rangle + \beta_t\eta_t^2\alpha^2\|\mathcal{A}\|^2. \end{aligned}$$

The second identity follows from the update rule (A.10). The bound on the last term depends on the fact that the constraint set $\alpha\Delta_n$ has Frobenius-norm diameter $\alpha\sqrt{2}$.

Next, we use the construction of the primal update to control the linear term in the last display. The update direction \mathbf{H}_t satisfies the inequality (A.9), so

$$(A.20) \quad \begin{aligned} \langle \mathbf{D}_t, \mathbf{H}_t - \mathbf{X}_t \rangle &\leq \min_{\mathbf{H} \in \alpha\Delta_n} \langle \mathbf{D}_t, \mathbf{H} - \mathbf{X}_t \rangle + \frac{\alpha\beta_0}{\beta_t}\|\mathbf{D}_t\| \\ &\leq \langle \mathbf{D}_t, \mathbf{X}_\star - \mathbf{X}_t \rangle + \frac{\alpha\beta_0}{\beta_t}\|\mathbf{D}_t\|. \end{aligned}$$

The second inequality depends on the fact that $\mathbf{X}_\star \in \alpha\Delta_n$.

We can use the explicit formula (A.8) for the derivative \mathbf{D}_t to control the two terms in (A.20). First,

$$(A.21) \quad \begin{aligned} \langle \mathbf{D}_t, \mathbf{X}_\star - \mathbf{X}_t \rangle &= \langle \mathbf{C} + \mathcal{A}^*(\mathbf{y}_t + \beta_t(\mathcal{A}\mathbf{X}_t - \mathbf{b})), \mathbf{X}_\star - \mathbf{X}_t \rangle \\ &= \langle \mathbf{C}, \mathbf{X}_\star - \mathbf{X}_t \rangle - \langle \mathbf{y}_t, \mathcal{A}\mathbf{X}_t - \mathbf{b} \rangle - \beta_t\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|^2 \\ &= \langle \mathbf{C}, \mathbf{X}_\star \rangle - f_t(\mathbf{X}_t) - \frac{1}{2}\beta_t\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|^2. \end{aligned}$$

We have invoked the definition of the adjoint \mathcal{A}^* and the fact that $\mathcal{A}\mathbf{X}_\star = \mathbf{b}$. Last, we used the definition (A.18) to identify the quantity $f_t(\mathbf{X}_t)$. Second,

$$(A.22) \quad \begin{aligned} \|\mathbf{D}_t\| &= \|\mathbf{C} + \mathcal{A}^*(\mathbf{y}_t + \beta_t(\mathcal{A}\mathbf{X}_t - \mathbf{b}))\| \\ &\leq \|\mathbf{C}\| + K\|\mathcal{A}\| + \beta_t\|\mathcal{A}\|\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|. \end{aligned}$$

We have used the assumption that \mathbf{y}_t satisfies the bounded travel condition (A.13).

Combine the last three displays to obtain the estimate

$$(A.23) \quad \begin{aligned} \langle \mathbf{D}_t, \mathbf{H}_t - \mathbf{X}_t \rangle &\leq \langle \mathbf{C}, \mathbf{X}_\star \rangle - f_t(\mathbf{X}_t) \\ &\quad + \left(\alpha\beta_0\|\mathcal{A}\|\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\| - \frac{1}{2}\beta_t\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|^2 \right) + \frac{\alpha\beta_0}{\beta_t}(\|\mathbf{C}\| + K\|\mathcal{A}\|). \end{aligned}$$

We can now control the decrease in the function f_t between adjacent primal iterates. Combine the displays (A.19) and (A.23) to arrive at the bound

$$\begin{aligned} f_t(\mathbf{X}_{t+1}) &\leq (1 - \eta_t)f_t(\mathbf{X}_t) + \eta_t\langle \mathbf{C}, \mathbf{X}_\star \rangle \\ &\quad + \left(\eta_t\alpha\beta_0\|\mathcal{A}\|\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\| - \frac{1}{2}\beta_t\eta_t\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|^2 \right) + \left(\beta_t\eta_t^2\alpha^2\|\mathcal{A}\|^2 + \frac{\eta_t\alpha\beta_0}{\beta_t}(\|\mathbf{C}\| + K\|\mathcal{A}\|) \right). \end{aligned}$$

Subtract $p_\star = \langle \mathbf{C}, \mathbf{X}_\star \rangle$ from both sides to arrive at

$$\begin{aligned} f_t(\mathbf{X}_{t+1}) - p_\star &\leq (1 - \eta_t)(f_t(\mathbf{X}_t) - p_\star) \\ &\quad + \left(\eta_t\alpha\beta_0\|\mathcal{A}\|\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\| - \frac{1}{2}\beta_t\eta_t\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|^2 \right) + \left(\beta_t\eta_t^2\alpha^2\|\mathcal{A}\|^2 + \frac{\eta_t\alpha\beta_0}{\beta_t}(\|\mathbf{C}\| + K\|\mathcal{A}\|) \right). \end{aligned}$$

Finally, use the definition (A.18) again to pass back to the augmented Lagrangian:

$$(A.24) \quad \begin{aligned} L_t(\mathbf{X}_{t+1}; \mathbf{y}_t) - p_\star &\leq (1 - \eta_t)(L_t(\mathbf{X}_t; \mathbf{y}_t) - p_\star) \\ &\quad + \left(\eta_t\alpha\beta_0\|\mathcal{A}\|\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\| - \frac{1}{2}\beta_t\eta_t\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|^2 \right) \\ &\quad + \left(\beta_t\eta_t^2\alpha^2\|\mathcal{A}\|^2 + \frac{\eta_t\alpha\beta_0}{\beta_t}(\|\mathbf{C}\| + K\|\mathcal{A}\|) \right). \end{aligned}$$

This bound describes the evolution of the augmented Lagrangian as the primal iterate advances. But we still need to include the effects of increasing the smoothing parameter and updating the dual variable.

A.6.2. Analysis of the smoothing update. To continue, observe that updates to the smoothing parameter have a controlled impact on the augmented Lagrangian (A.6):

$$L_t(\mathbf{X}_t; \mathbf{y}_t) - L_{t-1}(\mathbf{X}_t; \mathbf{y}_t) = \frac{1}{2}(\beta_t - \beta_{t-1})\|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|^2.$$

Add and subtract $L_{t-1}(\mathbf{X}_t; \mathbf{y}_t)$ in the large parenthesis in the first line of (A.24) and invoke the last identity to obtain

$$(A.25) \quad L_t(\mathbf{X}_{t+1}; \mathbf{y}_t) - p_\star \leq (1 - \eta_t) \left(L_{t-1}(\mathbf{X}_t; \mathbf{y}_t) - p_\star \right) \\ + \left(\eta_t \alpha \beta_0 \|\mathcal{A}\| \|\mathcal{A} \mathbf{X}_t - \mathbf{b}\| + \frac{1}{2} [(1 - \eta_t)(\beta_t - \beta_{t-1}) - \beta_t \eta_t] \|\mathcal{A} \mathbf{X}_t - \mathbf{b}\|^2 \right) \\ + \left(\beta_t \eta_t^2 \alpha^2 \|\mathcal{A}\|^2 + \frac{\eta_t \alpha \beta_0}{\beta_t} (\|\mathbf{C}\| + K \|\mathcal{A}\|) \right).$$

The next step is to develop a uniform bound on the terms in the second line so that we can ignore the role of the feasibility gap $\|\mathcal{A} \mathbf{X}_t - \mathbf{b}\|$ in the subsequent calculations. The choice (A.5) of the parameters ensures that

$$(1 - \eta_t)(\beta_t - \beta_{t-1}) - \beta_t \eta_t < \frac{-\beta_0^2}{\beta_t}.$$

Introduce this bound into the second line of (A.25) and maximize the resulting concave quadratic function to reach

$$\eta_t \alpha \beta_0 \|\mathcal{A}\| \|\mathcal{A} \mathbf{X}_t - \mathbf{b}\| + \frac{1}{2} [(1 - \eta_t)(\beta_t - \beta_{t-1}) - \beta_t \eta_t] \|\mathcal{A} \mathbf{X}_t - \mathbf{b}\|^2 \\ \leq \eta_t \alpha \beta_0 \|\mathcal{A}\| \|\mathcal{A} \mathbf{X}_t - \mathbf{b}\| - \frac{\beta_0^2}{2\beta_t} \|\mathcal{A} \mathbf{X}_t - \mathbf{b}\|^2 \leq \beta_t \eta_t^2 \alpha^2 \|\mathcal{A}\|^2.$$

Substitute the last display into (A.25) to determine that

$$(A.26) \quad L_t(\mathbf{X}_{t+1}; \mathbf{y}_t) - p_\star \leq (1 - \eta_t) \left(L_{t-1}(\mathbf{X}_t; \mathbf{y}_t) - p_\star \right) \\ + \left(2\beta_t \eta_t^2 \alpha^2 \|\mathcal{A}\|^2 + \frac{\eta_t \alpha \beta_0}{\beta_t} (\|\mathbf{C}\| + K \|\mathcal{A}\|) \right).$$

To develop a recursion, we need to assess how the left-hand side changes when we update the dual variable.

A.6.3. Analysis of the dual update. To incorporate the dual update, observe that

$$L_t(\mathbf{X}_{t+1}; \mathbf{y}_{t+1}) = L_t(\mathbf{X}_{t+1}; \mathbf{y}_t) + \langle \mathbf{y}_{t+1} - \mathbf{y}_t, \mathcal{A} \mathbf{X}_{t+1} - \mathbf{b} \rangle \\ = L_t(\mathbf{X}_{t+1}; \mathbf{y}_t) + \gamma_t \|\mathcal{A} \mathbf{X}_{t+1} - \mathbf{b}\|^2 \\ \leq L_t(\mathbf{X}_{t+1}; \mathbf{y}_t) + \beta_t \eta_t^2 \alpha^2 \|\mathcal{A}\|^2.$$

The first relation is simply the definition (A.6) of the augmented Lagrangian, while the second relation depends on the dual update rule (A.11). The last step follows from the selection rule (A.12) for the dual step size parameter.

Introduce the latter display into (A.25) to discover that

$$(A.27) \quad L_t(\mathbf{X}_{t+1}; \mathbf{y}_{t+1}) - p_\star \leq (1 - \eta_t) \left(L_{t-1}(\mathbf{X}_t; \mathbf{y}_t) - p_\star \right) \\ + \left(3\beta_t \eta_t^2 \alpha^2 \|\mathcal{A}\|^2 + \frac{\eta_t \alpha \beta_0}{\beta_t} (\|\mathbf{C}\| + K \|\mathcal{A}\|) \right).$$

We have developed a recursion for the value of the augmented Lagrangian as the iterates and the smoothing parameter evolve.

A.6.4. Solving the recursion. Next, we solve the recursion (A.27). We assert that

$$(A.28) \quad \begin{aligned} L_t(\mathbf{X}_{t+1}; \mathbf{y}_{t+1}) - p_\star &< \frac{2\beta_0}{\sqrt{t+1}} [6\alpha^2 \|\mathcal{A}\|^2 + \alpha(\|\mathbf{C}\| + K\|\mathcal{A}\|)] \\ &=: \frac{2\beta_0 E}{\sqrt{t+1}} \quad \text{for } t = 1, 2, 3, \dots \end{aligned}$$

For the case $t = 1$, the definition (A.5) ensures that $\eta_1 = 1$ and $\beta_1 = \beta_0\sqrt{2}$, so the bound (A.28) follows instantly from (A.27). When $t > 1$, an inductive argument using the recursion (A.27) and the bound (A.28) for $t - 1$ ensures that

$$\begin{aligned} L_t(\mathbf{X}_{t+1}; \mathbf{y}_{t+1}) - p_\star &\leq \left[\frac{t-1}{t+1} \cdot \frac{2\beta_0}{\sqrt{t}} + \frac{2\beta_0}{(t+1)^{3/2}} \right] [6\alpha^2 \|\mathcal{A}\|^2 + \alpha(\|\mathbf{C}\| + K\|\mathcal{A}\|)] \\ &< \frac{2\beta_0}{\sqrt{t+1}} [6\alpha^2 \|\mathcal{A}\|^2 + \alpha(\|\mathbf{C}\| + K\|\mathcal{A}\|)]. \end{aligned}$$

We have introduced the stated values (A.5) of the step size and smoothing parameters. The induction proceeds, and we conclude that (A.28) is valid.

A.6.5. Bound for the suboptimality of the objective. We are prepared to develop an upper bound on the suboptimality of the objective of the model problem (A.1). The definition (A.6) of the augmented Lagrangian directly implies that

$$(A.29) \quad \begin{aligned} \langle \mathbf{C}, \mathbf{X}_{t+1} \rangle - p_\star &= L_t(\mathbf{X}_{t+1}; \mathbf{y}_{t+1}) - p_\star - \langle \mathbf{y}_{t+1}, \mathcal{A}\mathbf{X}_{t+1} - \mathbf{b} \rangle \\ &\quad - \frac{1}{2}\beta_t \|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\|^2. \end{aligned}$$

Continuing from here,

$$\langle \mathbf{C}, \mathbf{X}_{t+1} \rangle - p_\star \leq \frac{2\beta_0 E}{\sqrt{t+1}} + K \cdot \|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\|.$$

The first identity follows from definition (A.6) of the augmented Lagrangian. The bound relies on (A.28) and the Cauchy–Schwarz inequality. We have also used the bounded travel condition (A.13). This establishes (A.14).

A.6.6. Bound for the superoptimality of the objective. The CGAL iterates \mathbf{X}_t are generally infeasible for (A.1), so they can be superoptimal. Nevertheless, we can easily control the superoptimality. By the saddle-point properties (A.3) and (A.4), the Lagrangian (A.2) satisfies

$$(A.30) \quad p_\star = \max_{\mathbf{y}} \min_{\mathbf{X} \in \alpha \Delta_n} L(\mathbf{X}; \mathbf{y}) \leq L(\mathbf{X}_{t+1}; \mathbf{y}_\star) = \langle \mathbf{C}, \mathbf{X}_{t+1} \rangle + \langle \mathbf{y}_\star, \mathcal{A}\mathbf{X}_{t+1} - \mathbf{b} \rangle.$$

Invoke the Cauchy–Schwarz inequality and rearrange to determine that

$$-\|\mathbf{y}_\star\| \cdot \|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\| \leq \langle \mathbf{C}, \mathbf{X}_{t+1} \rangle - p_\star.$$

This establishes (A.15).

A.6.7. Bound for the infeasibility of the iterates. Next, we demonstrate that the iterates converge toward the feasible set of (A.1). Combine (A.29) and (A.30) and rearrange to see that

$$\langle \mathbf{y}_{t+1} - \mathbf{y}_\star, \mathcal{A}\mathbf{X}_{t+1} - \mathbf{b} \rangle \leq L_t(\mathbf{X}_{t+1}; \mathbf{y}_{t+1}) - p_\star - \frac{1}{2}\beta_t \|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\|^2.$$

Bound the left-hand side below using Cauchy-Schwarz and the right-hand side above using (A.28):

$$-\|\mathbf{y}_{t+1} - \mathbf{y}_\star\| \cdot \|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\| \leq \frac{2\beta_0 E}{\sqrt{t+1}} - \frac{1}{2}\beta_t \|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\|^2.$$

Solve this quadratic inequality to obtain the bound

$$\begin{aligned} \|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\| &\leq \beta_t^{-1} \left(\|\mathbf{y}_{t+1} - \mathbf{y}_\star\| + \sqrt{\|\mathbf{y}_{t+1} - \mathbf{y}_\star\|^2 + \frac{4\beta_t\beta_0 E}{\sqrt{t+1}}} \right) \\ &\leq \beta_t^{-1} \left(2\|\mathbf{y}_{t+1} - \mathbf{y}_\star\| + \sqrt{4\beta_0^2 E} \right) \\ &= \frac{2\beta_0^{-1}\|\mathbf{y}_{t+1} - \mathbf{y}_\star\| + 2\sqrt{E}}{\sqrt{t+1}}. \end{aligned}$$

The second inequality depends on the definition (A.5) of the smoothing parameter β_t and the subadditivity of the square root.

Finally, we control the dual error using the bounded travel condition (A.13):

$$\|\mathbf{y}_{t+1} - \mathbf{y}_\star\| \leq \|\mathbf{y}_{t+1}\| + \|\mathbf{y}_\star\| \leq K + \|\mathbf{y}_\star\|.$$

The last two displays yield

$$\|\mathcal{A}\mathbf{X}_{t+1} - \mathbf{b}\| \leq \frac{2\beta_0^{-1}(K + \|\mathbf{y}_\star\|) + 2\sqrt{E}}{\sqrt{t+1}}.$$

This confirms (A.16).

A.7. Computable bounds for suboptimality. In this section, we assume that the linear minimization (A.9) is performed exactly at iteration t . Introduce the duality gap surrogate

$$(A.31) \quad g_t(\mathbf{X}) := \max_{\mathbf{H} \in \alpha \Delta_n} \langle \nabla f_t(\mathbf{X}), \mathbf{X} - \mathbf{H} \rangle.$$

The function f_t is defined in (A.18). Note that the gap $g(\mathbf{X}_t)$ can be evaluated with the information we have at hand:

$$(A.32) \quad \begin{aligned} g_t(\mathbf{X}_t) &= \langle \mathbf{D}_t, \mathbf{X}_t \rangle - \langle \mathbf{D}_t, \mathbf{H}_t \rangle \\ &= \langle \mathbf{C}, \mathbf{X}_t \rangle + \langle \mathbf{y}_t + \beta_t(\mathcal{A}\mathbf{X}_t - \mathbf{b}), \mathcal{A}\mathbf{X}_t \rangle - \langle \mathbf{D}_t, \mathbf{H}_t \rangle. \end{aligned}$$

The last term is just the value of the linear minimization (A.9)

The gap gives us computable bounds on the suboptimality of the current iterate \mathbf{X}_t . Indeed, the convexity of f_t implies that

$$f_t(\mathbf{X}_t) - f_t(\mathbf{X}_\star) \leq \langle \nabla f_t(\mathbf{X}_t), \mathbf{X}_t - \mathbf{X}_\star \rangle \leq g_t(\mathbf{X}_t).$$

Using the definition (A.18) and the fact that \mathbf{X}_\star is feasible for (A.1), we find that

$$f_t(\mathbf{X}_t) - p_\star \leq f_t(\mathbf{X}_t) - f_t(\mathbf{X}_\star) \leq g_t(\mathbf{X}_t).$$

Invoking the definition (A.18) again and rearranging, we find that

$$\begin{aligned} \langle \mathbf{C}, \mathbf{X}_t \rangle - p_\star &\leq g_t(\mathbf{X}_t) - \langle \mathbf{y}_t, \mathcal{A}\mathbf{X}_t - \mathbf{b} \rangle - \frac{1}{2}\beta_t \|\mathcal{A}\mathbf{X}_t - \mathbf{b}\|^2 \\ &\leq g_t(\mathbf{X}_t) - \langle \mathbf{y}_t, \mathcal{A}\mathbf{X}_t - \mathbf{b} \rangle. \end{aligned} \quad (\text{A.33})$$

In other words, the suboptimality of the primal iterate \mathbf{X}_t is bounded in terms of the dual gap $g_t(\mathbf{X}_t)$, the feasibility gap $\mathcal{A}\mathbf{X}_t - \mathbf{b}$, and the dual variable \mathbf{y}_t .

Appendix B. Sketching and reconstruction of a positive-semidefinite matrix.

This section reviews and gives additional details about the Nyström sketch [49, 46, 72, 100]. This sketch tracks an evolving psd matrix and then reports a provably accurate low-rank approximation. The material on error estimation is new.

B.1. Sketching and updates. Consider a psd input matrix $\mathbf{X} \in \mathbb{S}_n$. Let R be a parameter that modulates the storage cost of the sketch and the quality of the matrix approximation.

To construct the sketch, we draw and fix a standard normal test matrix $\mathbf{\Omega} \in \mathbb{F}^{n \times R}$. The sketch of the matrix \mathbf{X} takes the form

$$(\text{B.1}) \quad \mathbf{S} = \mathbf{X}\mathbf{\Omega} \in \mathbb{F}^{n \times R} \quad \text{and} \quad \tau = \text{tr}(\mathbf{X}) \in \mathbb{R}.$$

The sketch supports linear rank-one updates to \mathbf{X} . Indeed, we can track the evolution

$$\begin{aligned} \mathbf{X} &\leftarrow (1 - \eta) \mathbf{X} + \eta \mathbf{v} \mathbf{v}^* && \text{for } \eta \in [0, 1] \text{ and } \mathbf{v} \in \mathbb{F}^n \\ (\text{B.2}) \quad \text{via } \mathbf{S} &\leftarrow (1 - \eta) \mathbf{S} + \eta \mathbf{v} (\mathbf{v}^* \mathbf{\Omega}) \\ \text{and } \tau &\leftarrow (1 - \eta) \tau + \eta \|\mathbf{v}\|^2. \end{aligned}$$

The test matrix $\mathbf{\Omega}$ and the sketch (\mathbf{S}, τ) require storage of $2Rn + 1$ numbers in \mathbb{F} . The arithmetic cost of the linear update (5.2) to the sketch is $\Theta(Rn)$ numerical operations.

Remark B.1 (Structured random matrices). We can reduce storage costs by a factor of two by using a structured random matrix in place of $\mathbf{\Omega}$. For example, see [102, Sec. 3] or [98]. This modification requires additional care with implementation (e.g., use of sparse arithmetic or fast trigonometric transforms), but the improvement can be significant for very large problems.

B.2. The reconstruction process. Given the test matrix $\mathbf{\Omega}$ and the sketch $\mathbf{S} = \mathbf{X}\mathbf{\Omega}$, we can form a rank- R approximation $\widehat{\mathbf{X}}$ of the matrix \mathbf{X} contained in the sketch. This approximation is defined by the formula

$$(\text{B.3}) \quad \widehat{\mathbf{X}} := \mathbf{S}(\mathbf{\Omega}^* \mathbf{S})^\dagger \mathbf{S}^* = (\mathbf{X}\mathbf{\Omega})(\mathbf{\Omega}^* \mathbf{X}\mathbf{\Omega})^\dagger (\mathbf{X}\mathbf{\Omega})^*.$$

This reconstruction is called a *Nyström approximation*. We often truncate $\widehat{\mathbf{X}}$ by replacing it with its best rank- r approximation $\llbracket \widehat{\mathbf{X}} \rrbracket_r$ for some $r \leq R$.

See Algorithm B.1 for a numerically stable implementation of the Nyström reconstruction process (5.3), including error estimation steps. The algorithm takes $O(R^2n)$ numerical operations and $\Theta(Rn)$ storage.

Algorithm B.1 NystromSketch implementation (see [Appendix B](#))**Input:** Dimension n of input matrix, size R of sketch**Output:** Rank- R approximation $\widehat{\mathbf{X}}$ of sketched matrix in factored form $\widehat{\mathbf{X}} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$, where $\mathbf{U} \in \mathbb{R}^{n \times R}$ has orthonormal columns and $\mathbf{\Lambda} \in \mathbb{R}^{R \times R}$ is nonnegative diagonal, and the Schatten 1-norm approximation errors $\text{err}(\llbracket \widehat{\mathbf{X}} \rrbracket_r)$ for $1 \leq r \leq R$, as defined in [\(B.4\)](#)**Recommendation:** Choose R as large as possible

```

1 function NystromSketch.Init( $n, R$ )
2    $\mathbf{\Omega} \leftarrow \text{randn}(n, R)$  ▷ Draw and fix random test matrix
3    $\mathbf{S} \leftarrow \text{zeros}(n, R)$  and  $\tau \leftarrow 0$  ▷ Form sketch of zero matrix

4 function NystromSketch.RankOneUpdate( $\mathbf{v}, \eta$ ) ▷ Implements \(5.2\)
5    $\mathbf{S} \leftarrow (1 - \eta) \mathbf{S} + \eta \mathbf{v}(\mathbf{v}^* \mathbf{\Omega})$  ▷ Update sketch of matrix
6    $\tau \leftarrow (1 - \eta) \tau + \eta \|\mathbf{v}\|^2$  ▷ Update the trace

7 function NystromSketch.Reconstruct
8    $\sigma \leftarrow \sqrt{n} \text{eps}(\text{norm}(\mathbf{S}))$  ▷ Compute a shift parameter
9    $\mathbf{S} \leftarrow \mathbf{S} + \sigma \mathbf{\Omega}$  ▷ Implicitly form sketch of  $\mathbf{X} + \sigma \mathbf{I}$ 
10   $\mathbf{C} \leftarrow \text{chol}(\mathbf{\Omega}^* \mathbf{S})$ 
11   $(\mathbf{U}, \mathbf{\Sigma}, \sim) \leftarrow \text{svd}(\mathbf{S}/\mathbf{C})$  ▷ Dense SVD
12   $\mathbf{\Lambda} \leftarrow \max\{0, \mathbf{\Sigma}^2 - \sigma \mathbf{I}\}$  ▷ Remove shift
13   $\text{err} \leftarrow \tau - \text{cumsum}(\text{diag}(\mathbf{\Lambda}))$  ▷ Compute approximation errors

```

B.3. The approximation error. We can easily determine the exact Schatten 1-norm error in the truncated Nyström approximation:

$$(B.4) \quad \text{err}(\llbracket \widehat{\mathbf{X}} \rrbracket_r) := \|\mathbf{X} - \llbracket \widehat{\mathbf{X}} \rrbracket_r\|_* = \tau - \text{tr}(\llbracket \widehat{\mathbf{X}} \rrbracket_r) \quad \text{for each } r \leq R.$$

Furthermore, we can use [\(B.4\)](#) to ascertain whether the unknown input matrix \mathbf{X} is (almost) low-rank. Indeed, the best rank- r approximation of \mathbf{X} satisfies

$$(B.5) \quad \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_* \leq \text{err}(\llbracket \widehat{\mathbf{X}} \rrbracket_r) \quad \text{for each } r \leq R.$$

Thus, large drops in the function $r \mapsto \text{err}(\llbracket \widehat{\mathbf{X}} \rrbracket_r)$ signal large drops in the eigenvalues of \mathbf{X} . See [Appendix B.8](#) for the details.

B.4. Statistical properties of the Nyström sketch. The truncated Nyström approximation has a number of attractive statistical properties. For a fixed input matrix \mathbf{X} , the expected approximation error $\mathbb{E}_{\mathbf{\Omega}} \|\mathbf{X} - \llbracket \widehat{\mathbf{X}} \rrbracket_r\|$ is monotone decreasing in both the sketch size R and the truncation rank r . Furthermore, if $\text{rank}(\mathbf{X}) = r$ for $r \leq R$, then $\|\mathbf{X} - \llbracket \widehat{\mathbf{X}} \rrbracket_r\|_* = 0$ with probability one. We establish these results below in [Appendix B.9](#).

B.5. A priori error bounds. The Nyström approximation $\widehat{\mathbf{X}}$ yields a provably good estimate for the matrix \mathbf{X} contained in the sketch [[100](#), Thms. 4.1].

Fact B.2 (Nyström sketch: Error bound). Fix a psd matrix $\mathbf{X} \in \mathbb{S}_n$. Let $\mathbf{S} = \mathbf{X}\mathbf{\Omega}$ where $\mathbf{\Omega} \in \mathbb{F}^{n \times R}$ is standard normal. For each $r < R$, the Nyström approximation (B.3) satisfies

$$(B.6) \quad \mathbb{E}_{\mathbf{\Omega}} \|\mathbf{X} - \widehat{\mathbf{X}}\|_* \leq \left(1 + \frac{r}{R-r-1}\right) \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_*.$$

If we replace $\widehat{\mathbf{X}}$ with the rank- r truncation $\llbracket \widehat{\mathbf{X}} \rrbracket_r$, the error bound (B.6) remains valid. Similar results hold with high probability.

The truncated Nyström approximations satisfy a stronger error bound when the input matrix \mathbf{X} exhibits spectral decay.

Fact B.3 (Nyström sketch: Error bound II). Fix a psd matrix $\mathbf{X} \in \mathbb{S}_n$. Let $\mathbf{S} = \mathbf{X}\mathbf{\Omega}$ where $\mathbf{\Omega} \in \mathbb{F}^{n \times R}$ is standard normal. For each $r < R$, the Nyström approximation (B.3) satisfies

$$(B.7) \quad \mathbb{E}_{\mathbf{\Omega}} \|\mathbf{X} - \llbracket \widehat{\mathbf{X}} \rrbracket_r\|_* \leq \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_* + \left(1 + \frac{r}{R-r-1}\right) \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_*.$$

If we replace $\widehat{\mathbf{X}}$ with the rank- r truncation $\llbracket \widehat{\mathbf{X}} \rrbracket_r$, the error bound (B.7) remains valid. Similar results hold with high probability.

Proof. Combine the proofs of [100, Thm. 4.2] and [49, Thm. 9.3]. ■

B.6. Discussion. In practice, it is best to minimize the error attributable to sketching. To that end, we recommend choosing the sketch size parameter R as large as possible, given resource constraints, so that we can obtain the highest-quality Nyström approximation.

In some problems, e.g., MaxCut with eigenvector rounding, the desired rank r of the truncation is known in advance. In this case, Fact 5.2 offers guidance about how to select R to achieve a specific error tolerance $(1 + \zeta)$ in (2.3). For example, when $5r + 1 \leq R$, the expected Schatten 1-norm error in the rank- r approximation $\llbracket \widehat{\mathbf{X}} \rrbracket_r$ is at most $1.25 \times$ the error in the best rank- r approximation of \mathbf{X} .

When the input matrix \mathbf{X} has decaying eigenvalues, the error in the truncated approximation may be far smaller than Fact 5.2 predicts; see [100, Thm. 4.2]. This happy situation is typical when \mathbf{X} is generated by the CGAL iteration.

B.7. Representation of the truncated Nyström approximation. The key tool in the analysis is a simple representation for the truncated approximation. These facts are extracted from [100, Supp.].

Let $\mathbf{X} \in \mathbb{S}_n$ be a fixed psd matrix, and let $\mathbf{\Omega} \in \mathbb{R}^{n \times R}$ be an arbitrary test matrix. Let $\mathbf{P} \in \mathbb{S}_n$ be the orthoprojector onto the range of $\mathbf{X}^{1/2}\mathbf{\Omega}$. Then we can write the Nyström approximation (B.3) as

$$\widehat{\mathbf{X}} = \mathbf{X}^{1/2}\mathbf{P}\mathbf{X}^{1/2}.$$

For each $r \leq R$, define \mathbf{P}_r to be the orthoprojector onto the co-range of the matrix $\llbracket \mathbf{X}^{1/2}\mathbf{P} \rrbracket_r$. By construction, $\llbracket \mathbf{X}^{1/2}\mathbf{P} \rrbracket_r = \mathbf{X}^{1/2}\mathbf{P}_r$. As a consequence,

$$\llbracket \widehat{\mathbf{X}} \rrbracket_r = (\llbracket \mathbf{X}^{1/2}\mathbf{P} \rrbracket_r)(\llbracket \mathbf{X}^{1/2}\mathbf{P} \rrbracket_r)^* = \mathbf{X}^{1/2}\mathbf{P}_r\mathbf{X}^{1/2}.$$

These results allow us to relate the truncated approximations to each other:

B.8. Approximation errors: Analysis. We may now obtain explicit formulas for the error in each Nyström approximation. For $r \leq R$, note that

$$\mathbf{X} - \llbracket \widehat{\mathbf{X}} \rrbracket_r = \mathbf{X}^{1/2}(\mathbf{I} - \mathbf{P}_r)\mathbf{X}^{1/2} \succcurlyeq \mathbf{0}.$$

It follows immediately that

$$\text{err}(\llbracket \widehat{\mathbf{X}} \rrbracket_r) = \|\mathbf{X} - \llbracket \widehat{\mathbf{X}} \rrbracket_r\|_* = \text{tr}(\mathbf{X} - \llbracket \widehat{\mathbf{X}} \rrbracket_r) = \text{tr}(\mathbf{X}) - \text{tr}(\llbracket \widehat{\mathbf{X}} \rrbracket_r).$$

This is the relation (B.4).

Assume that $r \leq r'$. Since $\text{tr}(\llbracket \widehat{\mathbf{X}} \rrbracket_r) \leq \text{tr}(\llbracket \widehat{\mathbf{X}} \rrbracket_{r'})$, we have the bound

$$\text{err}(\llbracket \widehat{\mathbf{X}} \rrbracket_r) \geq \text{err}(\llbracket \widehat{\mathbf{X}} \rrbracket_{r'}) \quad \text{for } r \leq r'.$$

In other words, for fixed sketch size R , the error in the truncated Nyström approximation is monotone decreasing in the approximation rank.

To obtain (B.5), observe that

$$\|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_* \leq \|\mathbf{X} - \llbracket \widehat{\mathbf{X}} \rrbracket_r\|_* = \text{err}(\llbracket \widehat{\mathbf{X}} \rrbracket_r).$$

The inequality holds because $\llbracket \mathbf{X} \rrbracket_r$ is a best rank- r approximation of \mathbf{X} in Schatten 1-norm, while $\llbracket \widehat{\mathbf{X}} \rrbracket_r$ is another rank- r matrix.

B.9. Statistical properties: Analysis. Next, let us verify the statistical properties of the error. In this section, the test matrix $\boldsymbol{\Omega} \in \mathbb{F}^{n \times R}$ is standard normal.

Assuming that $\text{rank}(\mathbf{X}) = r \leq R$, let us prove that $\|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_* = 0$ with probability one. To that end, we observe

$$\text{range}(\mathbf{P}) = \text{range}(\mathbf{X}^{1/2}\boldsymbol{\Omega}) = \text{range}(\mathbf{X}^{1/2}) \quad \text{with probability one.}$$

It follows that

$$\widehat{\mathbf{X}} = \mathbf{X}^{1/2}\mathbf{P}\mathbf{X}^{1/2} = \mathbf{X} \quad \text{with probability one.}$$

Moreover, $\text{rank}(\mathbf{X}^{1/2}\mathbf{P}) = r$ with probability one. Conditional on this event,

$$\llbracket \widehat{\mathbf{X}} \rrbracket_r = \mathbf{X}^{1/2}\mathbf{P}_r\mathbf{X}^{1/2} = (\llbracket \mathbf{X}^{1/2}\mathbf{P} \rrbracket_r)(\llbracket \mathbf{X}^{1/2}\mathbf{P} \rrbracket_r)^* = (\mathbf{X}^{1/2}\mathbf{P})(\mathbf{X}^{1/2}\mathbf{P})^* = \mathbf{X}.$$

This is the stated result.

Next, we show that the expected error in the truncated Nyström approximation is monotone decreasing with respect to the sketch size R . Fix the truncation rank r . Let $\boldsymbol{\Omega}_+ = [\boldsymbol{\Omega} \quad \boldsymbol{\omega}]$, where $\boldsymbol{\omega} \in \mathbb{F}^n$ is a standard normal vector independent from $\boldsymbol{\Omega}$. Define $\mathbf{P}_+ \in \mathbb{S}_n$ to be the orthoprojector onto $\text{range}(\mathbf{X}^{1/2}\boldsymbol{\Omega}_+)$. It is clear that $\text{range}(\mathbf{P}) \subseteq \text{range}(\mathbf{P}_+)$, and so

$$\mathbf{X}^{1/2}\mathbf{P}\mathbf{X}^{1/2} \preceq \mathbf{X}^{1/2}\mathbf{P}_+\mathbf{X}^{1/2}.$$

As a consequence,

$$\text{tr}(\llbracket \mathbf{X}^{1/2}\mathbf{P}\mathbf{X}^{1/2} \rrbracket_r) \leq \text{tr}(\llbracket \mathbf{X}^{1/2}\mathbf{P}_+\mathbf{X}^{1/2} \rrbracket_r).$$

Equivalently,

$$\|X - \llbracket X^{1/2} P X^{1/2} \rrbracket_r\|_* \geq \|X - \llbracket X^{1/2} P_+ X^{1/2} \rrbracket_r\|_*.$$

Take the expectation with respect to Ω_+ . The left-hand side is the average error in the r -truncated Nyström approximation with a standard normal sketch of size R . The right-hand side is the same thing, except the sketch has size $R + 1$. This is the required result.

Appendix C. SketchyCGAL: Additional results. This section contains some additional material about the SketchyCGAL algorithm.

C.1. Solution quality. We can develop estimates for the quality of the SketchyCGAL solution by adapted the approach that we used for CGAL.

To do so, we need to track the primal objective value at the current iterate:

$$p_t = \langle C, X_t \rangle.$$

At each iteration, we can easily update this estimate using the computed approximate eigenvector v_t :

$$p_{t+1} = (1 - \eta_t)p_t + \eta_t \alpha \langle v_t, C v_t \rangle.$$

This update rule is applied with the help of the primitive (2.4)①.

When we wish to estimate the error, say in iteration t , we solve the eigenvalue subproblem to very high accuracy:

$$\xi_t = v_t^* D_t v_t = \min_{\|v\|=1} v^* D_t v.$$

Then, we can compute the surrogate duality gap:

$$g_t(X_t) = p_t + \langle y_t + \beta_t(z_t - b), z_t \rangle - \xi_t.$$

This expression follows directly from the formula (A.32) using the loop invariant that $z_t = \mathcal{A}X_t$. We arrive at a computable error estimate:

$$p_t - p_* \leq g_t(X_t) - \langle y_t, z_t - b \rangle.$$

This bound follows directly from (A.33).

C.2. Convergence theory. In this section, we establish two simple results on the convergence properties of the SketchyCGAL algorithm.

Theorem C.1 (SketchyCGAL: Convergence I). *Let Ψ_* be the solution set of the model problem (2.2). For each $r < R$, the iterates \widehat{X}_t computed by SketchyCGAL (subsections 6.2 and 6.3) satisfy*

$$\limsup_{t \rightarrow \infty} \mathbb{E}_\Omega \text{dist}_*(\widehat{X}_t, \Psi_*) \leq \left(1 + \frac{r}{R - r - 1}\right) \cdot \max_{X \in \Psi_*} \|X - \llbracket X \rrbracket_r\|_*.$$

The same bound holds for the truncated approximations $\llbracket \widehat{X}_t \rrbracket_r$. Here, dist_ is the nuclear-norm distance between a matrix and a set of matrices.*

Proof. The implicit iterates \mathbf{X}_t satisfy the conclusions of [Fact 3.1](#), so they converge toward the compact set Ψ_\star . Therefore, we can choose a sequence $\{\mathbf{X}_{t_\star}\} \subset \Psi_\star$ with the property that $\|\mathbf{X}_t - \mathbf{X}_{t_\star}\|_* \rightarrow 0$. By the triangle inequality and [\(6.9\)](#),

$$\begin{aligned} \mathbb{E}_\Omega \text{dist}_*(\widehat{\mathbf{X}}_t, \Psi_\star) &\leq \mathbb{E}_\Omega \|\widehat{\mathbf{X}}_t - \mathbf{X}_t\|_* + \text{dist}_*(\mathbf{X}_t, \Psi_\star) \\ &\leq \left(1 + \frac{r}{R-r-1}\right) \cdot \|\mathbf{X}_t - \llbracket \mathbf{X}_t \rrbracket_r\|_* + \|\mathbf{X}_t - \mathbf{X}_{t_\star}\|_*. \end{aligned}$$

The rank- r approximation error in Schatten 1-norm is 1-Lipschitz with respect to the Schatten 1-norm (cf. [\[102, Sec. SM2.2\]](#)), so

$$\|\mathbf{X}_t - \llbracket \mathbf{X}_t \rrbracket_r\|_* \leq \|\mathbf{X}_{t_\star} - \llbracket \mathbf{X}_{t_\star} \rrbracket_r\|_* + \|\mathbf{X}_t - \mathbf{X}_{t_\star}\|_* \leq \max_{\mathbf{X} \in \Psi_\star} \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_* + \|\mathbf{X}_t - \mathbf{X}_{t_\star}\|_*.$$

Combine the last two displays, and extract the superior limit. ■

If the implicit iterates generated by SketchyCGAL happen to converge to a limit, we have a more precise result.

Theorem C.2 (SketchyCGAL: Convergence II). *Assume the implicit iterates \mathbf{X}_t induced by SketchyCGAL ([subsections 6.2 and 6.3](#)) converge to a matrix \mathbf{X}_{cgal} . For each $r < R$, the computed iterates $\widehat{\mathbf{X}}_t$ satisfy*

$$\begin{aligned} \limsup_{t \rightarrow \infty} \mathbb{E}_\Omega \|\mathcal{A} \widehat{\mathbf{X}}_t - \mathbf{b}\| &\leq \left(1 + \frac{r}{R-r-1}\right) \cdot \|\mathcal{A}\| \cdot \|\mathbf{X}_{\text{cgal}} - \llbracket \mathbf{X}_{\text{cgal}} \rrbracket_r\|_*; \\ \limsup_{t \rightarrow \infty} \mathbb{E}_\Omega |\langle \mathbf{C}, \widehat{\mathbf{X}}_t \rangle - \langle \mathbf{C}, \mathbf{X}_\star \rangle| &\leq \left(1 + \frac{r}{R-r-1}\right) \cdot \|\mathbf{C}\| \cdot \|\mathbf{X}_{\text{cgal}} - \llbracket \mathbf{X}_{\text{cgal}} \rrbracket_r\|_*. \end{aligned}$$

The same bound holds for the truncated approximations $\llbracket \widehat{\mathbf{X}}_t \rrbracket_r$. If $\text{rank}(\mathbf{X}_{\text{cgal}}) \leq R$, then the computed iterates $\widehat{\mathbf{X}}_t$ converge to the solution set of [\(2.2\)](#).

Proof. The implicit iterates \mathbf{X}_t satisfy the conclusions of [Fact 3.1](#), so the limit \mathbf{X}_{cgal} solves [\(2.2\)](#). Using the triangle inequality, the operator norm bound, and [\(6.9\)](#), we obtain nonasymptotic error bounds

$$\begin{aligned} \mathbb{E}_\Omega \|\mathcal{A} \widehat{\mathbf{X}}_t - \mathbf{b}\| &\leq \frac{\text{Const}}{\sqrt{t}} + \left(1 + \frac{r}{R-r-1}\right) \cdot \|\mathcal{A}\| \cdot \|\mathbf{X}_t - \llbracket \mathbf{X}_t \rrbracket_r\|_*; \\ \mathbb{E}_\Omega |\langle \mathbf{C}, \widehat{\mathbf{X}}_t \rangle - \langle \mathbf{C}, \mathbf{X}_\star \rangle| &\leq \frac{\text{Const}}{\sqrt{t}} + \left(1 + \frac{r}{R-r-1}\right) \cdot \|\mathbf{C}\| \cdot \|\mathbf{X}_t - \llbracket \mathbf{X}_t \rrbracket_r\|_*. \end{aligned}$$

Extract the limit as $t \rightarrow \infty$. The last conclusion follows from the facts outlined in [Appendix B.5](#). ■

Appendix D. Beyond the model problem.

The CGAL algorithm [\[110\]](#) applies to a more general problem template than [\(2.2\)](#). Likewise, the SketchyCGAL algorithm can solve a wider class of problems in a scalable fashion. This section outlines some of the opportunities.

D.1. A more general template. Consider the optimization problem

$$(D.1) \quad \text{minimize } \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{subject to } \mathcal{A}\mathbf{X} \in \mathbf{K} \quad \text{and } \mathbf{X} \in \mathbf{X}, \quad \mathbf{X} \text{ is psd.}$$

In this expression, $\mathbf{K} \subset \mathbb{R}^d$ is a closed, convex set and $\mathbf{X} \subset \mathbb{S}_n(\mathbb{F})$ is a compact, convex set of matrices. The rest of this section describes some problems that fall within the compass of (D.1), as well as new computational challenges that appear. Algorithm D.1 contains pseudocode for a version of SketchyCGAL tailored to (D.1).

Remark D.1 (Other matrix optimization problems). We can also extend SketchyCGAL to optimization problems involving matrices that are symmetric (but not psd) or that are rectangular. For example, matrix completion via nuclear-norm minimization [96] falls in this framework. In this case, we need to replace the Nyström sketch with a more general technique, such as [101, 102]. Further extensions are also possible; see [110]. We omit these developments.

D.2. The convex constraint set. To handle the convex constraint \mathbf{X} that appears in (D.1), we must develop a subroutine for the linear minimization problem

$$(D.2) \quad \underset{\mathbf{H} \in \mathbb{S}_n}{\text{minimize}} \quad \langle \mathbf{D}_t, \mathbf{H} \rangle \quad \text{subject to } \mathbf{H} \in \mathbf{X}, \quad \mathbf{H} \text{ is psd.}$$

To implement SketchyCGAL efficiently, we need the problem (D.2) to admit a structured (e.g., low-rank) approximate solution. Here are some situations where this is possible.

1. **Trace-bounded psd matrices.** $\mathbf{X} := \{\mathbf{X} \in \mathbb{S}_n : \text{tr } \mathbf{X} \leq \alpha \text{ and } \mathbf{X} \text{ is psd}\}$. For solving a standard-form SDP, this constraint is more natural than $\mathbf{X} = \alpha \mathbf{\Delta}_n$. Given an (approximate) minimum eigenpair (ξ_t, \mathbf{v}_t) of \mathbf{D}_t , the solution of (D.2) is

$$\mathbf{H}_t = \begin{cases} \alpha \mathbf{v}_t \mathbf{v}_t^*, & \xi_t < 0, \\ \mathbf{0}, & \xi_t \geq 0. \end{cases}$$

As before, we can solve the eigenvector problem with Algorithms 4.1 and 4.2.

2. **Relaxed orthoprojectors.** $\mathbf{X} := \{\mathbf{X} \in \mathbb{S}_n : \text{tr } \mathbf{X} = \alpha \text{ and } \mathbf{0} \preceq \mathbf{X} \preceq \mathbf{I}\}$. This is the best convex relaxation of the set of orthogonal projectors with rank α ; see [85]. When α is small, we can provably solve the linear minimization with randomized subspace iteration [49] or randomized block Lanczos methods [48, Sec. 10.3.6].

D.3. Convex inclusions. To handle the inclusion in \mathbf{K} that appears in (D.1), we need an efficient algorithm to perform the Euclidean projection onto \mathbf{K} . That is,

$$\text{proj}_{\mathbf{K}}(\mathbf{w}) := \arg \min \{\|\mathbf{w} - \mathbf{u}\| : \mathbf{u} \in \mathbf{K}\} \quad \text{for } \mathbf{w} \in \mathbb{R}^d.$$

Here are some important examples:

1. **Inequality constraints.** $\mathbf{K} := \{\mathbf{u} \in \mathbb{R}^d : \mathbf{u} \leq \mathbf{b}\}$. In this case, the projection takes the form $\text{proj}_{\mathbf{K}}(\mathbf{w}) = (\mathbf{w} - \mathbf{b})_-$, where $(\cdot)_-$ reports the negative part of a vector.
2. **Norm constraints.** $\mathbf{K} := \{\mathbf{u} \in \mathbb{R}^d : \|\mathbf{u} - \mathbf{b}\| \leq \delta\}$, where $\|\cdot\|$ is a norm. The projector can be computed easily for many norms, including the ℓ_p norm for $p \in \{1, 2, \infty\}$.

D.4. The CGAL iteration for the general template. To extend the description of the CGAL iteration in [Appendix A.3](#) for the general template [\(D.1\)](#), we consider the following augmented Lagrangian formulation with the slack variable $\mathbf{w} \in \mathbf{K}$ instead of [\(A.6\)](#):

$$L_t(\mathbf{X}; \mathbf{y}) := \langle \mathbf{C}, \mathbf{X} \rangle + \min_{\mathbf{w} \in \mathbf{K}} \left\{ \langle \mathbf{y}, \mathcal{A}\mathbf{X} - \mathbf{w} \rangle + \frac{1}{2} \beta_t \|\mathcal{A}\mathbf{X} - \mathbf{w}\|^2 \right\}.$$

Accordingly, the partial derivative [\(A.8\)](#) becomes

$$\mathbf{D}_t := \partial_{\mathbf{X}} L_t(\mathbf{X}_t; \mathbf{y}_t) = \mathbf{C} + \mathcal{A}^*(\mathbf{y}_t + \beta_t(\mathcal{A}\mathbf{X}_t - \mathbf{w}_t)) \quad \text{where} \quad \mathbf{w}_t := \text{proj}_{\mathbf{K}}(\mathcal{A}\mathbf{X}_t + \beta_t^{-1}\mathbf{y}_t).$$

We replace the linear minimization subroutine [\(A.9\)](#) with [\(D.2\)](#). We can still use an inexact variant of [\(D.2\)](#) with additive error. We also modify the dual update scheme by modifying [\(A.11\)](#) as

$$\mathbf{y}_{t+1} = \mathbf{y}_t + \gamma_t(\mathcal{A}\mathbf{X}_{t+1} - \bar{\mathbf{w}}_t) \quad \text{where} \quad \bar{\mathbf{w}}_t := \text{proj}_{\mathbf{K}}(\mathcal{A}\mathbf{X}_{t+1} + \beta_{t+1}^{-1}\mathbf{y}_t).$$

Finally, we replace the dual step size parameter selection rule [\(A.12\)](#) with

$$(D.3) \quad \gamma_t \|\mathcal{A}\mathbf{X}_{t+1} - \bar{\mathbf{w}}_t\|^2 \leq \beta_t \eta_t^2 \alpha^2 \|\mathcal{A}\|^2.$$

The bounded travel condition [\(A.13\)](#) remains the same.

To obtain the extension of SketchyCGAL to the general template, we simply pursue the same program outlined in [section 6](#) to augment CGAL with sketching.

Appendix E. Details of phase retrieval experiments. This section presents further details about the phase retrieval experiments presented in [section 7](#).

E.1. Synthetic phase retrieval data. This section provides additional details on the construction of synthetic datasets for the abstract phase retrieval SDP.

For each $n \in \{10^2, 10^3, \dots, 10^6\}$, we generate 20 independent datasets as follows. First, draw $\mathbf{x}_{\mathfrak{h}} \in \mathbb{C}^n$ from the complex standard normal distribution. We acquire $d = 12n$ phaseless measurements [\(7.2\)](#) using the coded diffraction pattern model [\[31\]](#).

To do so, we randomly draw 12 independent modulating waveforms ψ_j for $j = 1, 2, \dots, 12$. Each entry of ψ_j is drawn as the product of two independent random variables, one chosen uniformly from $\{1, i, -1, -i\}$, and the other from $\{\sqrt{2}/2, \sqrt{3}\}$ with probabilities 0.8 and 0.2 respectively. Then, we modulate $\mathbf{x}_{\mathfrak{h}}$ with these waveforms and take its Fourier transform. Each \mathbf{a}_i corresponds to computing a single entry of this Fourier transform:

$$\mathbf{a}_{(j-1)n+\ell} = \mathbf{W}_n(\ell, :) \text{diag}^*(\psi_j), \quad 1 \leq j \leq 12 \quad \text{and} \quad 1 \leq \ell \leq n,$$

where $\mathbf{W}_n(\ell, :)$ is the ℓ th row of the $n \times n$ discrete Fourier transform matrix. We use the fast Fourier transform to implement the measurement operator.

E.2. Fourier Ptychography. We study a more realistic measurement setup, Fourier ptychography (FP), for the phase retrieval problem. In this setup, $\mathbf{x}_{\mathfrak{h}} \in \mathbb{C}^n$ corresponds to an unknown high resolution image (vectorized) from a microscopic sample in the Fourier domain.

Algorithm D.1 SketchyCGAL for the general template (D.1)

Input: Problem data for (D.1) implemented via the primitives (2.4), sketch size R , number T of iterations

Output: Rank- R approximate solution to (D.1) in factored form $\widehat{\mathbf{X}}_T = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^*$ where $\mathbf{U} \in \mathbb{R}^{n \times R}$ has orthonormal columns and $\mathbf{\Lambda} \in \mathbb{R}^{R \times R}$ is nonnegative diagonal, and the Schatten 1-norm approximation errors $\text{err}(\|\widehat{\mathbf{X}}\|_r)$ for $1 \leq r \leq R$, as defined in (B.4)

Recommendation: To achieve (2.3), set R as large as possible, and set $T \approx \varepsilon^{-1}$

```

1 function SketchyCGAL( $R; T$ )
2   Scale problem data (subsection 7.1.2)                                ▷ [opt] Recommended!
3    $\beta_0 \leftarrow 1$  and  $K \leftarrow +\infty$                                 ▷ Default parameters
4   NystromSketch.Init( $n, R$ )
5    $\mathbf{z} \leftarrow \mathbf{0}_d$  and  $\mathbf{y} \leftarrow \mathbf{0}_d$ 
6   for  $t \leftarrow 1, 2, 3, \dots, T$  do
7      $\beta \leftarrow \beta_0 \sqrt{t+1}$  and  $\eta \leftarrow 2/(t+1)$ 
8      $\mathbf{w} \leftarrow \text{proj}_K(\mathbf{z} + \beta^{-1}\mathbf{y})$ 
9      $\mathbf{D} \leftarrow \mathbf{C} + \mathcal{A}^*(\mathbf{y} + \beta(\mathbf{z} - \mathbf{w}))$                                 ▷ Represent via primitives (2.4)❶❷
10     $\mathbf{H}$  is a (low-rank) matrix that solves (D.2)
11     $\mathbf{z} \leftarrow (1 - \eta)\mathbf{z} + \eta\mathcal{A}(\mathbf{H})$                                 ▷ Use primitive (2.4)❸
12     $\beta_+ \leftarrow \beta_0 \sqrt{t+2}$ 
13     $\mathbf{w} \leftarrow \text{proj}_K(\mathbf{z} + \beta_+^{-1}\mathbf{y})$ 
14     $\mathbf{y} \leftarrow \mathbf{y} + \gamma(\mathbf{z} - \mathbf{w})$                                 ▷ Step size  $\gamma$  satisfies (A.13) and (D.3)
15    NystromSketch.RankOneUpdate( $\sqrt{\alpha}\mathbf{v}, \eta$ )
16  ( $\mathbf{U}, \mathbf{\Lambda}, \text{err}$ )  $\leftarrow$  NystromSketch.Reconstruct

```

One cannot directly acquire a high resolution image from this sample because of the physical limitations of optical systems. Any measurement is subject to a filter caused by the lens aperture. We can represent this filter by a sparse matrix $\Phi \in \mathbb{C}^{m \times n}$ with $m \leq n$, each row of which has only one non-zero coefficient. Because of this filter, we can acquire only low-resolution images, through m -dimensional Fourier transform.

FP enlightens the sample from L different angles using a LED grid. This lets us to obtain L different aperture matrices Φ_j , $j = 1, 2, \dots, L$. Then, we acquire phaseless measurements from the sample using the following transmission matrices:

$$\mathbf{a}_{(j-1)m+\ell} = \mathbf{W}_m^*(\ell, :) \Phi_j, \quad 1 \leq j \leq L \quad \text{and} \quad 1 \leq \ell \leq m.$$

Here, $\mathbf{W}_m^*(\ell, :)$ is the ℓ th row of the conjugate transpose of discrete Fourier transform matrix.

The aim in FP is to reconstruct complex valued $\chi_{\mathfrak{t}}$ from these phaseless measurements. Once we construct $\chi_{\mathfrak{t}}$, we can generate a high resolution image by taking its inverse Fourier transform.

Appendix F. Additional numerical results.

This section contains quantitative results from the MaxCut and QAP experiments.

F.1. The MaxCut SDP. This section gives further information about the MaxCut experiments summarized in subsection 7.2. Table 3 presents numerical data from the MaxCut SDP experiment with GSET benchmark. We compare the methods in terms of the cut weight, objective value, primal and dual infeasibility, and the storage cost and computation time. The storage cost is approximated by monitoring the virtual memory size of the process, hence it includes the memory that is swapped out and it can go beyond 16 GB.

F.2. Failure of the Burer–Monteiro heuristic. This section presents empirical evidence that Burer–Monteiro factorization methods cannot support storage costs better than $\Omega(n\sqrt{d})$. Our approach is based on the paper of Waldspurger & Waters [105], which proves that the BM heuristic (8.2) can produce incorrect results unless $R = \Omega(n\sqrt{d})$.

Waldspurger provided us code that generates a random symmetric $\mathbf{C} \in \mathbb{S}_n(\mathbb{R})$. The MaxCut SDP (1.3) with objective \mathbf{C} has a unique solution, and the solution has rank 1. If the factorization rank R satisfies $R(R+3) \leq 2n$, then the BM formulation (8.2) has second-order critical points that are not optimal points of the original SDP (1.3). As a consequence, the Burer–Monteiro approach is reliable only if the storage budget is $\Theta(n^{3/2})$. In contrast, for the same problem instances, our analysis (Theorem Theorem 6.3) shows that SketchyCGAL succeeds with factorization rank $R = 2$ and storage budget $\Theta(n)$.

We will demonstrate numerically that Waldspurger & Waters [105] have identified a serious obstruction to using the Burer–Monteiro approach. Moreover, we will see that SketchyCGAL resolves the issue. See the code supplement for scripts to reproduce these experiments.

We use the Manopt software [22] to solve the Burer–Monteiro formulation of the MaxCut SDP. For $n = 100$, we drew 10 random matrices $\mathbf{C}_1, \dots, \mathbf{C}_{10}$ using Waldspurger’s code. For each instance, we sweep the factorization rank $R = 2, 3, 4, \dots, 13$. (For $R \geq 13$, we anticipate that each second-order critical point of the Burer–Monteiro problem is a solution to the original SDP, owing to the analysis in [23].) In each experiment, we ran Manopt with 100 random initializations, and we counted the number of times the algorithm failed. We declared failure if Manopt converged to a second-order stationary point whose objective value is 10^{-3} larger than the true optimal value. See Table 2 for the statistics.

In contrast, SketchyCGAL can solve all of these instances, even when the sketch size $R = 2$. For these problems, we use the default parameter choices for SketchyCGAL, but we do not pre-scale the data or perform tuning. Figure F.1 compares the convergence trajectory of SketchyCGAL and Manopt for one problem instance. The difference is evident.

F.3. The quadratic assignment problem. This section gives further information about the QAP experiments summarized in subsection 7.5. Tables 5 and 7 display the performance of SketchyCGAL, by presenting the upper bound after rounding (subsection 7.5.3), the objective value $\langle \mathbf{B} \otimes \mathbf{A}, \mathbf{X} \rangle$, the feasibility gap $\text{dist}_K(\mathcal{A}\mathbf{X})$, the total number of iterations, the memory usage, and the computation time.

Tables 4 and 6 compare the relative gap (7.7) obtained by SketchyCGAL with the values for the CSDP method [26] with clique size $k = \{2, 3, 4\}$ and the PATH method [116] reported in [26, Tab. 6]. These results appeared also in Figure 7.5.

Acknowledgments. We would like to thank Irène Waldspurger for providing code that generates instances of the MaxCut SDP that are difficult for the Burer–Monteiro heuristic.

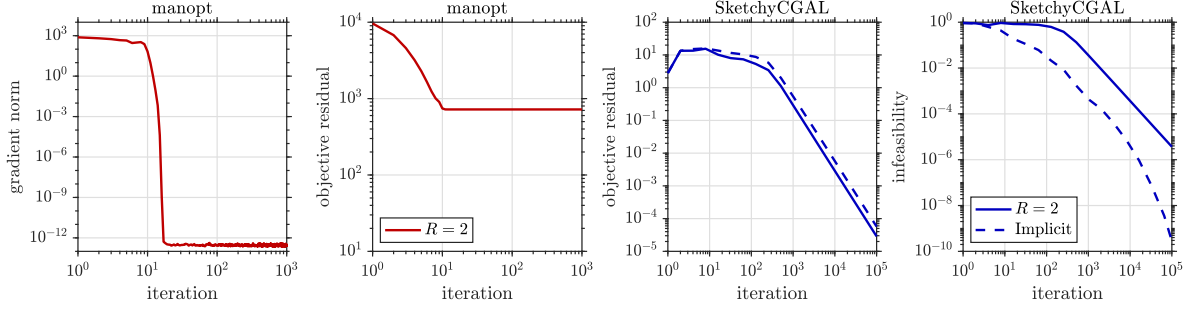


Figure F.1. MaxCut SDP: Failure of the Burer–Monteiro heuristic. We apply *Manopt* and *SketchyCGAL* for solving the MaxCut SDP with the dataset C_1 . The subplots show the gradient norm and suboptimality for *Manopt* [left] and the suboptimality and infeasibility for *SketchyCGAL* [right]. *Manopt* with $R = 2$ converges to a spurious solution, whereas *SketchyCGAL* successfully computes a rank-1 approximation of the global optimum. The dashed line describes the convergence of the *SketchyCGAL* implicit iterates. For details, see [Appendix F.2](#).

Table 2

We run *Manopt* for solving hard instances of the MaxCut SDP. We consider 10 datasets C_1, \dots, C_{10} . For each dataset, we run *Manopt* with 100 random initializations and report the number of failures. We declare failure when *Manopt* converges to second-order critical point that is not a global optimum. For details, see [Appendix F.2](#).

Dataset / R	R = 2	3	4	5	6	7	8	9	10	11	12	13
C_1	82	69	63	53	35	32	24	12	11	1	4	0
C_2	77	56	56	36	19	17	12	2	0	0	0	0
C_3	89	65	54	47	44	46	23	11	5	0	3	0
C_4	84	69	50	40	27	23	18	17	1	0	9	0
C_5	85	68	52	51	43	30	31	20	14	3	4	0
C_6	81	68	53	41	23	22	10	10	2	0	1	0
C_7	83	76	60	39	19	19	19	3	0	0	1	0
C_8	81	73	44	34	41	25	8	12	5	4	10	0
C_9	84	64	46	35	25	17	1	10	0	2	4	0
C_{10}	83	71	54	50	31	25	24	16	13	0	8	0

Table 3: Numerical outcomes from the MaxCut SDP experiment with GSET Benchmark.

	SketchyCGAL	MoSeK	SDPT3	SDPNAL+	Sedumi
Data Name Size (n)	cut weight $\langle \mathbf{C}, \mathbf{X} \rangle$ $\ \mathcal{A}\mathbf{X} - \mathbf{b}\ / \max\{\ \mathbf{b}\ , 1\}$ $ \min\{\lambda_{\min}(\mathbf{X}), 0\} $ storage (MB) cpu time (hh:mm:ss)				
G1 800	11382 -1.2072e+04 2.7955e-03 0 7 00:00:00:03	11414 -1.2083e+04 3.0760e-08 0 225 00:00:00:15	11414 -1.2083e+04 9.3301e-13 0 180 00:00:00:15	11414 -1.2083e+04 8.6224e-15 0 258 00:00:01:35	11414 -1.2083e+04 3.6902e-09 0 354 00:00:00:43
G2 800	11361 -1.2082e+04 2.8124e-03 0 7 00:00:00:03	11329 -1.2089e+04 3.4843e-08 0 225 00:00:00:16	11329 -1.2089e+04 9.5745e-13 0 180 00:00:00:12	11339 -1.2089e+04 3.3489e-15 0 275 00:00:01:43	11329 -1.2089e+04 2.9720e-09 0 352 00:00:00:41
G3 800	11388 -1.2076e+04 2.5103e-03 0 7 00:00:00:03	11392 -1.2084e+04 5.9077e-08 0 225 00:00:00:19	11392 -1.2084e+04 1.7491e-11 0 180 00:00:00:13	11392 -1.2084e+04 3.4555e-07 0 219 00:00:01:15	11392 -1.2084e+04 3.9665e-09 0 356 00:00:00:51
G4 800	11387 -1.2103e+04 2.6474e-03 0 7 00:00:00:05	11366 -1.2111e+04 4.2044e-08 0 225 00:00:00:14	11366 -1.2111e+04 1.8551e-13 0 180 00:00:00:11	11368 -1.2111e+04 3.3027e-07 0 277 00:00:01:16	11366 -1.2111e+04 3.3830e-09 0 354 00:00:00:52
G5 800	11356 -1.2092e+04 2.6573e-03 0 7 00:00:00:03	11432 -1.2100e+04 4.2128e-08 0 225 00:00:00:14	11432 -1.2100e+04 3.1188e-13 0 180 00:00:00:12	11432 -1.2100e+04 6.3811e-09 0 270 00:00:01:50	11432 -1.2100e+04 3.2630e-09 0 352 00:00:00:45
G6 800	1899 -2.6534e+03 2.7480e-03 0 7 00:00:00:06	1950 -2.6562e+03 2.8773e-08 0 225 00:00:00:15	1950 -2.6562e+03 3.2757e-12 0 180 00:00:00:16	1950 -2.6562e+03 1.0055e-07 0 273 00:00:01:17	1950 -2.6562e+03 1.7342e-09 0 354 00:00:00:43
G7 800	1779 -2.4873e+03 2.7858e-03 0 7 00:00:00:04	1749 -2.4893e+03 2.3310e-08 0 225 00:00:00:22	1749 -2.4893e+03 2.5344e-12 0 180 00:00:00:16	1749 -2.4893e+03 1.2441e-07 0 268 00:00:01:14	1749 -2.4893e+03 1.6999e-09 0 352 00:00:01:06

Continue on the next page

Table 3: MaxCut Benchmark with GSet (cont.).

	SketchyCGAL	MoSeK	SDPT3	SDPNAL+	Sedumi
G8 800	1784	1765	1765	1765	1765
	-2.5056e+03	-2.5069e+03	-2.5069e+03	-2.5069e+03	-2.5069e+03
	2.7816e-03	4.2417e-08	3.4138e-13	8.5771e-07	1.7960e-09
	0	0	0	0	0
	7	225	180	277	289
	00:00:00:04	00:00:00:15	00:00:00:15	00:00:01:25	00:00:00:45
G9 800	1776	1800	1800	1800	1800
	-2.5262e+03	-2.5287e+03	-2.5287e+03	-2.5287e+03	-2.5287e+03
	2.7310e-03	2.9469e-08	8.2962e-12	2.4830e-07	3.9707e-10
	0	0	0	0	0
	7	225	180	277	352
	00:00:00:06	00:00:00:14	00:00:00:12	00:00:01:08	00:00:00:46
G10 800	1790	1766	1766	1766	1766
	-2.4840e+03	-2.4851e+03	-2.4851e+03	-2.4851e+03	-2.4851e+03
	2.7810e-03	3.7738e-08	5.7206e-13	7.4784e-07	2.9698e-10
	0	0	0	0	0
	7	225	180	277	352
	00:00:00:04	00:00:00:15	00:00:00:12	00:00:01:21	00:00:01:04
G11 800	524	512	512	522	512
	-6.2814e+02	-6.2916e+02	-6.2916e+02	-6.2914e+02	-6.2916e+02
	2.6719e-03	1.7634e-13	1.4007e-12	1.0682e-10	1.8847e-09
	0	0	0	0	0
	7	225	186	277	352
	00:00:00:04	00:00:00:19	00:00:00:09	00:00:01:52	00:00:00:39
G12 800	524	512	512	514	512
	-6.2279e+02	-6.2387e+02	-6.2387e+02	-6.2387e+02	-6.2387e+02
	2.6153e-03	1.8468e-13	5.3073e-13	1.8087e-11	1.1260e-09
	0	0	0	0	0
	7	225	186	278	352
	00:00:00:03	00:00:00:14	00:00:00:07	00:00:02:20	00:00:00:46
G13 800	554	534	534	534	534
	-6.4594e+02	-6.4714e+02	-6.4714e+02	-6.4713e+02	-6.4714e+02
	2.6652e-03	1.6420e-13	8.3034e-12	3.4895e-15	7.7471e-10
	0	0	0	0	0
	7	225	186	278	354
	00:00:00:03	00:00:00:14	00:00:00:08	00:00:01:47	00:00:00:41
G14 800	2953	2967	2967	2963	2967
	-3.1821e+03	-3.1916e+03	-3.1916e+03	-3.1916e+03	-3.1916e+03
	2.4329e-03	3.4503e-08	4.1785e-13	7.3214e-09	1.1476e-09
	0	0	0	0	0
	7	225	180	284	352
	00:00:00:05	00:00:00:16	00:00:00:09	00:00:01:47	00:00:00:49
G15 800	2948	2971	2971	2971	2971
	-3.1611e+03	-3.1716e+03	-3.1716e+03	-3.1716e+03	-3.1716e+03
	2.1250e-03	5.9551e-08	9.9206e-14	3.8762e-07	1.6534e-09
	0	0	0	0	0
	7	225	186	278	352
	00:00:00:04	00:00:00:15	00:00:00:13	00:00:02:23	00:00:00:45

Continue on the next page

Table 3: MaxCut Benchmark with GSet (cont.).

	SketchyCGAL	MoSeK	SDPT3	SDPNAL+	Sedumi
G16 800	2967	2958	2958	2956	2958
	-3.1658e+03	-3.1750e+03	-3.1750e+03	-3.1750e+03	-3.1750e+03
	2.3354e-03	8.1927e-08	4.7805e-13	9.3060e-15	3.1805e-10
	0	0	0	0	0
	7	225	180	279	352
	00:00:00:04	00:00:00:20	00:00:00:10	00:00:02:00	00:00:00:46
G17 800	2938	2957	2957	2960	2957
	-3.1621e+03	-3.1713e+03	-3.1713e+03	-3.1713e+03	-3.1713e+03
	2.4028e-03	1.3790e-08	1.2021e-12	9.2606e-07	7.2377e-10
	0	0	0	0	0
	7	225	186	280	352
	00:00:00:04	00:00:00:17	00:00:00:10	00:00:02:17	00:00:00:49
G18 800	898	893	893	893	893
	-1.1643e+03	-1.1660e+03	-1.1660e+03	-1.1660e+03	-1.1660e+03
	2.8264e-03	1.0695e-08	1.2377e-11	9.8235e-07	1.2709e-09
	0	0	0	0	0
	7	225	186	278	352
	00:00:00:06	00:00:00:17	00:00:00:09	00:00:01:37	00:00:01:06
G19 800	792	797	797	797	797
	-1.0790e+03	-1.0820e+03	-1.0820e+03	-1.0820e+03	-1.0820e+03
	2.7808e-03	1.4813e-08	5.5509e-12	3.8521e-08	1.3036e-09
	0	0	0	0	0
	7	225	186	273	352
	00:00:00:06	00:00:00:16	00:00:00:12	00:00:01:38	00:00:00:58
G20 800	846	837	837	837	837
	-1.1086e+03	-1.1114e+03	-1.1114e+03	-1.1114e+03	-1.1114e+03
	2.7607e-03	1.6650e-08	8.3457e-13	6.6975e-07	7.4903e-10
	0	0	0	0	0
	7	225	186	285	354
	00:00:00:06	00:00:00:17	00:00:00:10	00:00:01:34	00:00:01:02
G21 800	811	841	841	841	841
	-1.1019e+03	-1.1043e+03	-1.1043e+03	-1.1043e+03	-1.1043e+03
	2.6359e-03	9.8852e-09	5.5307e-13	4.8156e-07	2.0021e-10
	0	0	0	0	0
	7	225	186	268	357
	00:00:00:05	00:00:00:17	00:00:00:13	00:00:01:33	00:00:00:54
G22 2000	12974	12956	12956	12955	12956
	-1.4125e+04	-1.4136e+04	-1.4136e+04	-1.4136e+04	-1.4136e+04
	4.3829e-03	1.0335e-08	7.0113e-13	1.9644e-07	2.7849e-09
	0	0	0	0	0
	8	1049	735	994	1498
	00:00:00:09	00:00:06:17	00:00:01:35	00:00:49:50	00:00:17:38
G23 2000	12918	12987	12987	12996	12987
	-1.4133e+04	-1.4142e+04	-1.4142e+04	-1.4142e+04	-1.4142e+04
	4.3010e-03	1.6980e-08	1.6925e-14	1.6469e-09	8.3706e-10
	0	0	0	0	0
	8	1049	705	865	1401
	00:00:00:04	00:00:04:12	00:00:02:36	00:00:46:53	00:00:14:02

Continue on the next page

Table 3: MaxCut Benchmark with GSet (cont.).

	SketchyCGAL	MoSeK	SDPT3	SDPNAL+	Sedumi
G24 2000	12885	12979	12979	12978	12979
	-1.4135e+04	-1.4141e+04	-1.4141e+04	-1.4141e+04	-1.4141e+04
	4.4016e-03	1.0315e-08	1.1767e-12	2.8584e-09	7.7807e-10
	0	0	0	0	0
	8	1049	670	962	1498
	00:00:00:10	00:00:04:07	00:00:02:03	00:00:46:51	00:00:14:26
G25 2000	12913	12885	12885	12897	12885
	-1.4139e+04	-1.4144e+04	-1.4144e+04	-1.4144e+04	-1.4144e+04
	4.3776e-03	1.7277e-08	5.8895e-13	7.4190e-08	1.1674e-09
	0	0	0	0	0
	8	1049	705	993	1401
	00:00:00:17	00:00:06:31	00:00:01:31	00:00:46:09	00:00:16:08
G26 2000	12847	12866	12866	12866	12866
	-1.4122e+04	-1.4133e+04	-1.4133e+04	-1.4133e+04	-1.4133e+04
	4.3789e-03	1.0479e-08	5.0789e-14	4.6882e-08	1.0125e-09
	0	0	0	0	0
	8	1049	705	993	1466
	00:00:00:06	00:00:03:57	00:00:01:48	00:00:47:04	00:00:26:36
G27 2000	2863	2888	2888	2888	2888
	-4.1378e+03	-4.1417e+03	-4.1417e+03	-4.1417e+03	-4.1417e+03
	4.4638e-03	1.5899e-08	5.9893e-12	2.5844e-07	1.2743e-09
	0	0	0	0	0
	8	1049	768	1001	1466
	00:00:00:19	00:00:04:12	00:00:03:03	00:00:27:23	00:00:29:58
G28 2000	2856	2843	2843	2843	2843
	-4.0971e+03	-4.1008e+03	-4.1008e+03	-4.1008e+03	-4.1008e+03
	4.4429e-03	9.6101e-09	3.6450e-12	8.0676e-08	1.0742e-09
	0	0	0	0	0
	8	1049	734	1001	1432
	00:00:00:14	00:00:04:10	00:00:01:45	00:00:50:35	00:00:16:56
G29 2000	2999	2978	2978	2978	2978
	-4.2056e+03	-4.2089e+03	-4.2089e+03	-4.2089e+03	-4.2089e+03
	4.4590e-03	4.7152e-08	6.4617e-12	2.6237e-07	1.2095e-09
	0	0	0	0	0
	8	1049	705	973	1465
	00:00:00:13	00:00:04:32	00:00:01:49	00:00:46:07	00:00:26:32
G30 2000	2993	3007	3007	3007	3007
	-4.2107e+03	-4.2154e+03	-4.2154e+03	-4.2154e+03	-4.2154e+03
	4.3487e-03	2.8211e-08	1.1262e-11	3.6633e-07	1.3249e-09
	0	0	0	0	0
	8	1049	714	1033	1465
	00:00:00:07	00:00:07:03	00:00:02:54	00:00:29:38	00:00:27:07
G31 2000	2840	2876	2876	2876	2876
	-4.1139e+03	-4.1167e+03	-4.1167e+03	-4.1167e+03	-4.1167e+03
	4.4544e-03	1.3367e-08	4.0493e-11	9.9629e-07	1.4877e-09
	0	0	0	0	0
	8	1049	768	1040	1432
	00:00:00:09	00:00:04:01	00:00:02:27	00:00:52:58	00:00:20:18

Continue on the next page

Table 3: MaxCut Benchmark with GSet (cont.).

	SketchyCGAL	MoSeK	SDPT3	SDPNAL+	Sedumi
G32 2000	1286	1268	1268	1288	1268
	-1.5645e+03	-1.5676e+03	-1.5676e+03	-1.5675e+03	-1.5676e+03
	4.4160e-03	2.1823e-13	3.5478e-11	2.9972e-14	8.8110e-10
	0	0	0	0	0
	8	1049	589	974	1466
	00:00:00:05	00:00:03:57	00:00:00:58	00:01:13:59	00:00:23:45
G33 2000	1262	1276	1276	1282	1276
	-1.5410e+03	-1.5443e+03	-1.5443e+03	-1.5442e+03	-1.5443e+03
	4.4238e-03	1.8467e-13	2.1554e-11	5.0275e-15	2.4434e-09
	0	0	0	0	0
	8	1049	595	1036	1432
	00:00:00:05	00:00:03:36	00:00:01:02	00:01:20:22	00:00:12:44
G34 2000	1266	1250	1254	1258	1256
	-1.5435e+03	-1.5467e+03	-1.5467e+03	-1.5466e+03	-1.5467e+03
	4.4154e-03	2.9281e-13	1.1258e-12	9.1706e-15	3.4482e-10
	0	0	0	0	0
	8	1049	595	1005	1466
	00:00:00:05	00:00:04:06	00:00:00:56	00:01:26:33	00:00:13:43
G35 2000	7418	7411	7411	7411	7411
	-8.0057e+03	-8.0147e+03	-8.0147e+03	-8.0147e+03	-8.0147e+03
	4.4657e-03	2.3327e-08	1.3377e-11	5.6863e-09	6.6740e-10
	0	0	0	0	0
	8	1049	654	993	1494
	00:00:00:18	00:00:04:36	00:00:01:46	00:00:57:44	00:00:21:25
G36 2000	7371	7399	7399	7396	7399
	-7.9941e+03	-8.0060e+03	-8.0060e+03	-8.0060e+03	-8.0060e+03
	4.4700e-03	3.5444e-08	2.6518e-14	2.2431e-09	4.4800e-10
	0	0	0	0	0
	8	1049	590	994	1498
	00:00:00:17	00:00:06:47	00:00:01:55	00:00:58:12	00:00:16:56
G37 2000	7374	7443	7443	7437	7443
	-8.0007e+03	-8.0186e+03	-8.0186e+03	-8.0186e+03	-8.0186e+03
	4.4645e-03	3.7915e-08	8.6663e-12	7.5503e-09	2.0343e-09
	0	0	0	0	0
	8	1049	724	993	1401
	00:00:00:19	00:00:04:24	00:00:01:21	00:00:48:17	00:00:18:56
G38 2000	7370	7402	7402	7400	7402
	-7.9999e+03	-8.0150e+03	-8.0150e+03	-8.0150e+03	-8.0150e+03
	4.4400e-03	2.7416e-08	4.4579e-13	8.8096e-09	2.0761e-10
	0	0	0	0	0
	8	1049	590	888	1431
	00:00:00:19	00:00:06:27	00:00:02:25	00:00:48:12	00:00:19:18
G39 2000	2143	2176	2176	2171	2176
	-2.8700e+03	-2.8776e+03	-2.8776e+03	-2.8776e+03	-2.8776e+03
	4.4570e-03	1.0983e-08	1.7938e-12	1.7759e-09	5.1775e-10
	0	0	0	0	0
	8	1049	621	1067	1563
	00:00:00:15	00:00:04:37	00:00:01:49	00:00:47:07	00:00:32:23

Continue on the next page

Table 3: MaxCut Benchmark with GSet (cont.).

	SketchyCGAL	MoSeK	SDPT3	SDPNAL+	Sedumi
G40 2000	2071	2153	2153	2151	2153
	-2.8611e+03	-2.8648e+03	-2.8648e+03	-2.8648e+03	-2.8648e+03
	4.4469e-03	3.9814e-08	2.8335e-11	3.3046e-07	1.6695e-09
	0	0	0	0	0
	8	1049	723	1067	1532
	00:00:00:22	00:00:05:05	00:00:01:31	00:00:29:13	00:00:32:19
G41 2000	2087	2114	2114	2114	2114
	-2.8602e+03	-2.8652e+03	-2.8652e+03	-2.8652e+03	-2.8652e+03
	4.4599e-03	4.6842e-08	1.5034e-12	6.0544e-09	7.0174e-10
	0	0	0	0	0
	8	1049	641	1067	1498
	00:00:00:15	00:00:06:12	00:00:01:31	00:00:50:39	00:00:31:37
G42 2000	2186	2172	2172	2178	2172
	-2.9392e+03	-2.9463e+03	-2.9463e+03	-2.9463e+03	-2.9463e+03
	4.4565e-03	1.7392e-08	4.3777e-12	1.8720e-15	1.1033e-09
	0	0	0	0	0
	8	1049	671	936	1468
	00:00:00:12	00:00:06:28	00:00:01:28	00:00:25:19	00:00:21:00
G43 1000	6524	6518	6518	6518	6518
	-7.0268e+03	-7.0322e+03	-7.0322e+03	-7.0322e+03	-7.0322e+03
	3.0606e-03	4.4803e-08	1.1908e-13	6.6168e-07	1.3780e-09
	0	0	0	0	0
	7	298	246	347	483
	00:00:00:02	00:00:00:28	00:00:00:18	00:00:02:23	00:00:02:08
G44 1000	6491	6461	6461	6461	6461
	-7.0241e+03	-7.0279e+03	-7.0279e+03	-7.0279e+03	-7.0279e+03
	3.0788e-03	7.0255e-09	1.0885e-11	9.6232e-07	4.9394e-10
	0	0	0	0	0
	7	298	246	340	420
	00:00:00:02	00:00:00:28	00:00:00:22	00:00:02:26	00:00:01:28
G45 1000	6449	6447	6447	6447	6447
	-7.0208e+03	-7.0248e+03	-7.0248e+03	-7.0248e+03	-7.0248e+03
	3.0408e-03	1.0257e-08	9.2987e-12	9.7518e-07	6.2000e-10
	0	0	0	0	0
	7	298	246	348	420
	00:00:00:04	00:00:00:41	00:00:00:20	00:00:03:02	00:00:01:29
G46 1000	6458	6416	6416	6416	6416
	-7.0250e+03	-7.0299e+03	-7.0299e+03	-7.0299e+03	-7.0299e+03
	3.1575e-03	1.4069e-08	7.8291e-13	7.3592e-07	1.6385e-09
	0	0	0	0	0
	7	298	246	344	483
	00:00:00:04	00:00:00:32	00:00:00:18	00:00:02:23	00:00:01:45
G47 1000	6461	6452	6452	6454	6452
	-7.0324e+03	-7.0367e+03	-7.0367e+03	-7.0367e+03	-7.0367e+03
	3.0751e-03	4.6358e-09	5.9188e-12	6.8425e-15	1.1436e-09
	0	0	0	0	0
	7	298	246	348	482
	00:00:00:02	00:00:00:29	00:00:00:18	00:00:03:01	00:00:01:38

Continue on the next page

Table 3: MaxCut Benchmark with GSet (cont.).

	SketchyCGAL	MoSeK	SDPT3	SDPNAL+	Sedumi
G48 3000	6000	6000	6000	6000	6000
	-5.9870e+03	-6.0000e+03	-6.0000e+03	-6.0812e+03	-6.0000e+03
	5.2439e-03	2.1863e-14	2.5288e-13	2.1707e-13	2.0368e-12
	0	0	0	0	0
	9	2086	959	1390	2621
	00:00:00:03	00:00:07:24	00:00:02:13	00:02:52:47	00:00:32:41
G49 3000	6000	6000	6000	6000	6000
	-5.9871e+03	-6.0000e+03	-6.0000e+03	-6.0026e+03	-6.0000e+03
	5.4493e-03	1.4375e-13	3.3753e-14	3.9549e-13	4.4770e-10
	0	0	0	0	0
	9	2086	959	1424	2621
	00:00:00:02	00:00:07:44	00:00:01:49	00:01:53:02	00:00:43:07
G50 3000	5868	5880	5880	5880	5880
	-5.9796e+03	-5.9882e+03	-5.9882e+03	-5.9915e+03	-5.9882e+03
	5.4017e-03	5.4964e-14	3.4115e-14	9.6276e-14	2.0697e-13
	0	0	0	0	0
	9	2152	959	1400	2621
	00:00:00:02	00:00:08:45	00:00:02:20	00:02:46:35	00:00:29:42
G51 1000	3687	3738	3738	3737	3738
	-3.9918e+03	-4.0063e+03	-4.0063e+03	-4.0063e+03	-4.0063e+03
	2.2432e-03	6.5908e-08	4.7342e-13	9.3109e-07	1.4928e-09
	0	0	0	0	0
	7	298	246	348	483
	00:00:00:06	00:00:00:28	00:00:00:18	00:00:03:40	00:00:02:29
G52 1000	3727	3711	3711	3708	3711
	-3.9956e+03	-4.0096e+03	-4.0096e+03	-4.0096e+03	-4.0096e+03
	2.3113e-03	1.1607e-08	1.1565e-12	1.0138e-06	2.1532e-09
	0	0	0	0	0
	7	298	246	347	418
	00:00:00:06	00:00:00:30	00:00:00:23	00:00:03:51	00:00:01:27
G53 1000	3714	3712	3712	3712	3712
	-3.9972e+03	-4.0097e+03	-4.0097e+03	-4.0097e+03	-4.0097e+03
	1.9289e-03	6.2655e-08	6.4444e-14	1.0682e-14	4.4861e-10
	0	0	0	0	0
	7	298	246	355	483
	00:00:00:06	00:00:00:28	00:00:00:17	00:00:03:31	00:00:02:11
G54 1000	3700	3716	3716	3718	3716
	-3.9934e+03	-4.0062e+03	-4.0062e+03	-4.0062e+03	-4.0062e+03
	1.9950e-03	1.4837e-08	1.4319e-12	1.5096e-14	1.5897e-09
	0	0	0	0	0
	7	298	311	347	482
	00:00:00:07	00:00:00:43	00:00:00:14	00:00:03:35	00:00:01:47
G55 5000	9884	9866	9866	9865	9866
	-1.1036e+04	-1.1039e+04	-1.1039e+04	-1.1039e+04	-1.1039e+04
	4.0738e-03	6.5842e-09	1.0881e-12	2.1875e-07	1.9760e-09
	0	0	0	0	0
	11	6541	2646	3701	7448
	00:00:02:42	00:01:08:53	00:00:13:54	00:08:41:41	00:04:30:44

Continue on the next page

Table 3: MaxCut Benchmark with GSet (cont.).

	SketchyCGAL	MoSeK	SDPT3	SDPNAL+	Sedumi
G56 5000	3584	3595	3595	3596	3595
	-4.7569e+03	-4.7600e+03	-4.7600e+03	-4.7600e+03	-4.7600e+03
	7.0679e-03	3.2421e-08	4.2710e-14	1.1317e-11	7.1008e-10
	0	0	0	0	0
	11	6541	2735	3700	6862
	00:00:01:49	00:01:05:47	00:00:19:58	00:13:19:39	00:04:12:45
G57 5000	3150	3152	3156	3192	3158
	-3.8791e+03	-3.8855e+03	-3.8855e+03	-3.8853e+03	-3.8855e+03
	7.0600e-03	1.6506e-13	2.7378e-13	4.6605e-14	1.7408e-09
	0	0	0	0	0
	11	6541	2652	3700	6861
	00:00:00:16	00:00:55:50	00:00:13:39	00:20:51:45	00:03:49:29
G58 5000	18470	18561	18561	18559	18561
	-2.0135e+04	-2.0136e+04	-2.0136e+04	-2.0136e+04	-2.0136e+04
	7.0648e-03	1.4230e-08	2.4310e-12	4.5090e-15	9.3470e-10
	0	0	0	0	0
	11	6541	2748	3724	6862
	00:00:02:52	00:01:17:25	00:00:18:26	00:14:15:52	00:04:40:20
G59 5000	5321	5325	5325	5316	5325
	-7.2866e+03	-7.3123e+03	-7.3123e+03	-7.3123e+03	-7.3123e+03
	7.0688e-03	1.0705e-08	4.5632e-12	4.1272e-15	1.3622e-09
	0	0	0	0	0
	11	6541	2802	3734	6862
	00:00:01:46	00:01:14:29	00:00:19:02	00:18:10:33	00:09:19:12
G60 7000	13617	13610	13610	13623	13610
	-1.5219e+04	-1.5222e+04	-1.5222e+04	-1.5222e+04	-1.5222e+04
	7.0468e-03	1.0901e-08	1.1320e-12	5.4916e-14	8.9220e-10
	0	0	0	0	0
	20	12681	5017	7074	13237
	00:00:04:04	00:02:27:21	00:00:53:31	01:10:39:19	00:11:31:27
G61 7000	5208	5218	5218	5194	5218
	-6.8239e+03	-6.8281e+03	-6.8281e+03	-6.8281e+03	-6.8281e+03
	8.3648e-03	1.7210e-08	5.3273e-13	2.8846e-15	1.3178e-09
	0	0	0	0	0
	20	12681	5395	7137	13237
	00:00:02:09	00:02:37:09	00:00:41:19	01:08:33:11	00:17:32:16
G62 7000	4406	4384	4384	4404	4384
	-5.4258e+03	-5.4309e+03	-5.4309e+03	-5.4306e+03	-5.4309e+03
	8.3568e-03	1.6462e-13	6.5021e-13	3.2675e-14	2.3616e-09
	0	0	0	0	0
	16	12681	5089	7071	13236
	00:00:00:56	00:02:04:21	00:00:35:16	01:13:44:46	00:09:23:00
G63 7000	25890	25988	25988	26000	25988
	-2.8218e+04	-2.8244e+04	-2.8244e+04	-2.8244e+04	-2.8244e+04
	8.3635e-03	4.4392e-08	1.6210e-11	6.7982e-15	1.4974e-09
	0	0	0	0	0
	16	12681	5235	7077	13237
	00:00:06:04	00:02:48:34	00:00:52:03	01:17:16:18	00:12:58:10

Continue on the next page

Table 3: MaxCut Benchmark with GSet (cont.).

	SketchyCGAL	MoSeK	SDPT3	SDPNAL+	Sedumi
G64 7000	7645	7746	7746	7745	7746
	-1.0438e+04	-1.0466e+04	-1.0466e+04	-1.0466e+04	-1.0466e+04
	8.3622e-03	4.9100e-08	2.5557e-11	1.7486e-14	1.5548e-09
	0	0	0	0	0
	20	12681	5649	7076	13237
	00:00:02:10	00:02:51:13	00:00:50:25	01:14:32:37	00:14:01:50
G65 8000	5062	5008	5010	5026	—
	-6.1967e+03	-6.2055e+03	-6.2055e+03	-6.2051e+03	—
	8.9344e-03	4.0923e-13	2.4743e-11	1.7440e-14	—
	0	0	0	0	—
	16	15495	6613	9177	—
	00:00:00:36	00:03:37:10	00:00:52:36	02:17:39:55	—
G66 9000	5750	5764	5766	5782	—
	-7.0624e+03	-7.0772e+03	-7.0772e+03	-7.0768e+03	—
	9.4618e-03	6.4942e-13	2.3288e-12	1.4708e-14	—
	0	0	0	0	—
	18	21969	8405	11565	—
	00:00:00:31	00:04:51:33	00:01:14:52	03:07:36:48	—
G67 10000	6266	—	6224	6278	—
	-7.7367e+03	—	-7.7444e+03	-7.7441e+03	—
	9.9992e-03	—	1.6629e-14	1.6371e-08	—
	0	—	0	0	—
	16	—	10335	14231	—
	00:00:02:13	—	00:01:38:47	05:21:16:12	—

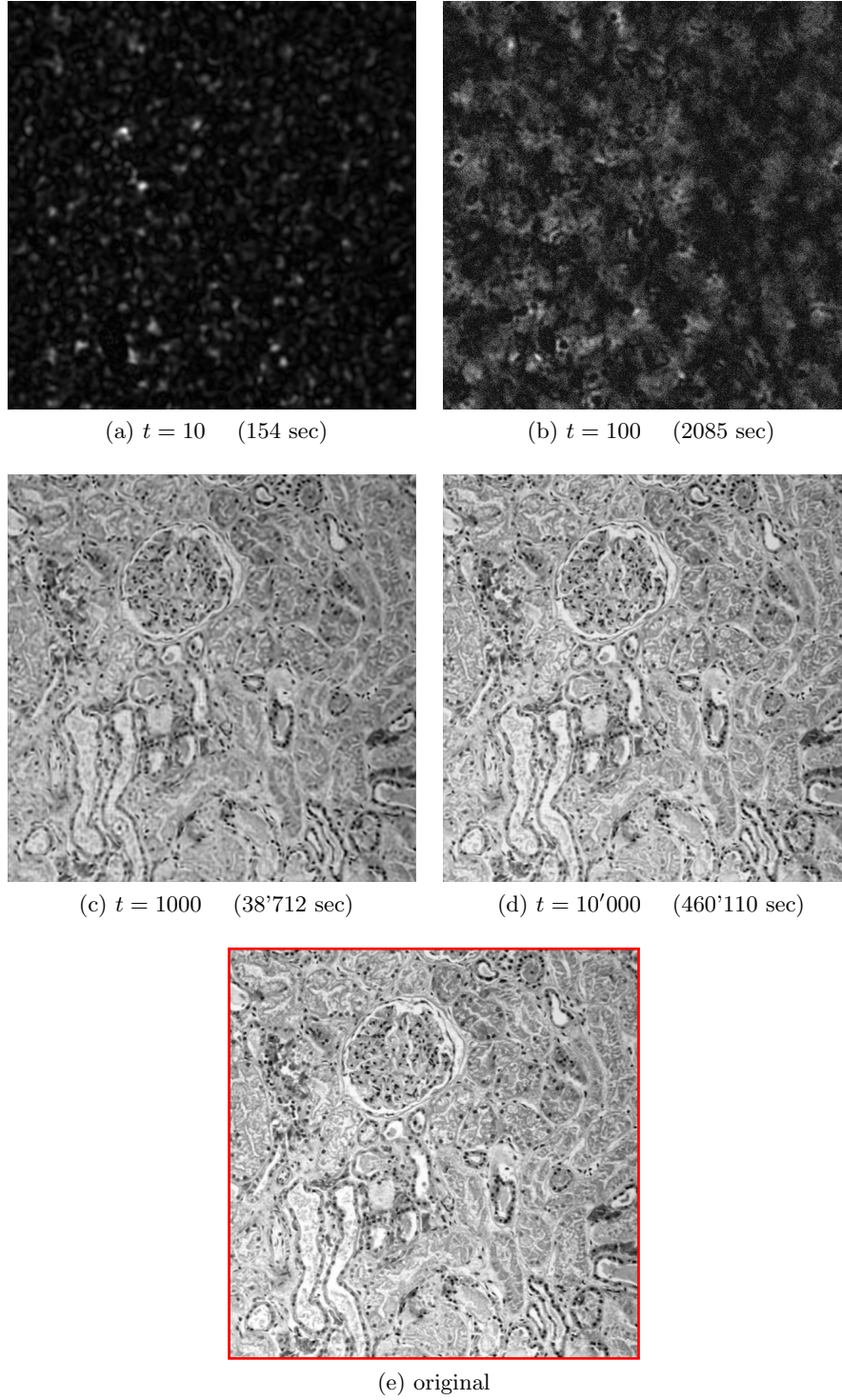


Figure F.2. Phase retrieval SDP: Imaging. Reconstruction of an $n = 640^2$ pixel image from Fourier ptychography data. We solve an $n \times n$ phase retrieval SDP via *SketchyCGAL* with rank parameter $R = 5$ and show the images obtained at iterations $t = 10, 10^2, 10^3, 10^4$. The last subfigure is the original. See [subsection 7.4.3](#).

Table 4

We solve SDP relaxations of QAP instances from QAPLIB using SketchyCGAL. We compute the relative gap and compare it with the values for the CSDP method [26] with clique size $k = \{2, 3, 4\}$ and the PATH method [116] reported in [26, Tab. 4]. Smaller is better. See subsection 7.5.

Dataset	Optimum	SketchyCGAL	CSDP2	CSDP3	CSDP4	PATH
chr12a	9552	0	34.5	6	0	42.7
chr12b	9742	0	38.9	25.4	11.9	38.1
chr12c	11156	0	5.8	2.3	2.3	18.6
chr15a	9896	0.4	2.1	2.1	2.1	52
chr15b	7990	0	26.3	34.5	29.2	158.6
chr15c	9504	0	0	0	0	63.3
chr18a	11098	1.5	69.8	0.2	0.2	76.3
chr18b	1534	8.7	8.9	22.9	29.5	99.3
chr20a	2192	6.6	122.5	76.1	43.8	95.4
chr20b	2298	0	62.9	9.3	9.3	82.2
chr20c	14142	17.1	173	100.1	111.5	88.9
chr22a	6156	0	17.2	7.6	3	38.3
chr22b	6194	4.7	7.3	2.3	1	40.4
chr25a	3796	21.8	107	49.2	25.2	69.9
esc16a	68	2.9	8.8	11.8	11.8	11.8
esc16b	292	0	0	0.7	0	2.7
esc16c	160	3.8	5	7.5	8.7	6.3
esc16d	16	12.5	12.5	50	25	75
esc16e	28	14.3	14.3	7.1	14.3	21.4
esc16g	26	0	7.7	0	15.4	15.4
esc16h	996	0	1.6	0	1.6	16.9
esc16i	14	0	0	0	0	57.1
esc16j	8	0	0	0	0	75
esc32a	130	52.3	115.4	124.6	113.8	93.8
esc32b	168	64.3	109.5	114.3	111.9	88.1
esc32c	642	2.2	12.8	15.9	13.7	7.8
esc32d	200	14	38	36	39	21
esc32e	2	0	0	0	0	600
esc32g	6	0	0	0	0	366.7
esc32h	438	5	24.7	26.9	22.8	18.3
esc64a	116	6.9	60.3	53.4	60.3	106.9
esc128	64	28.1	250	206.3	175	221.9
ste36a	9526	15.1	70.2	74.7	74.2	76.3
ste36b	15852	15.7	188.8	204.3	211.9	158.6
ste36c	8239110	9.7	66	62.8	63.7	83.2

Table 5

We run *SketchyCGAL* for (the first of) 10^6 iterations or 48 hours of runtime, for solving SDP relaxations of QAP instances from QAPLIB. We report the upper bound after rounding (subsection 7.5.3), the objective value, the feasibility gap, the number of iterations, the memory usage (in MB), and the cpu time ('hh:mm:ss'). See subsection 7.5.

Dataset	Optimum	Upper bnd.	Objective	Feas. gap	Iteration	Memory	Time
chr12a	9552	9552	9 576,06	5.02e−5	1000000	69	02 : 58 : 40
chr12b	9742	9742	9 759,43	2.43e−5	1000000	69	02 : 57 : 36
chr12c	11156	11156	11 021,43	2.05e−4	1000000	69	02 : 58 : 09
chr15a	9896	9936	9 519,70	1.91e−4	1000000	82	04 : 43 : 27
chr15b	7990	7990	7 469,99	1.72e−4	1000000	82	04 : 43 : 35
chr15c	9504	9504	9 517,70	3.70e−5	1000000	82	04 : 43 : 29
chr18a	11098	11262	10 700,38	2.39e−4	1000000	82	07 : 05 : 36
chr18b	1534	1668	1 534,30	8.08e−5	1000000	82	07 : 42 : 28
chr20a	2192	2336	2 183,67	1.11e−4	1000000	69	10 : 10 : 16
chr20b	2298	2298	2 289,02	9.43e−5	1000000	69	10 : 12 : 48
chr20c	14142	16554	13 316,80	1.89e−4	1000000	69	10 : 14 : 18
chr22a	6156	6156	6 152,78	1.06e−4	1000000	65	13 : 01 : 58
chr22b	6194	6486	6 198,40	9.91e−5	1000000	65	13 : 05 : 36
chr25a	3796	4624	3 692,54	2.06e−4	1000000	65	24 : 57 : 56
esc16a	68	70	60,36	5.36e−5	1000000	69	07 : 52 : 29
esc16b	292	292	287,64	9.63e−4	1000000	65	12 : 22 : 31
esc16c	160	166	145,01	5.59e−5	1000000	65	09 : 11 : 47
esc16d	16	18	13,00	6.37e−5	1000000	69	06 : 12 : 40
esc16e	28	32	25,41	6.89e−5	1000000	69	06 : 15 : 12
esc16g	26	26	22,46	7.05e−5	1000000	69	06 : 23 : 49
esc16h	996	996	975,41	2.06e−4	1000000	65	12 : 34 : 09
esc16i	14	14	11,37	8.25e−5	1000000	69	05 : 45 : 60
esc16j	8	8	7,11	8.49e−5	1000000	69	05 : 22 : 29
esc32a	130	198	99,60	3.65e−4	799533	208	48 : 00 : 00
esc32b	168	276	118,36	2.12e−4	634813	130	48 : 00 : 00
esc32c	642	656	610,85	7.33e−4	496050	65	48 : 00 : 00
esc32d	200	228	187,33	2.95e−4	609878	130	48 : 00 : 00
esc32e	2	2	1,90	4.67e−4	1000000	65	32 : 39 : 09
esc32g	6	6	5,83	4.48e−4	1000000	65	35 : 43 : 20
esc32h	438	460	417,78	3.79e−4	468726	65	48 : 00 : 00
esc64a	116	124	97,78	7.00e−3	163634	162	48 : 00 : 00
esc128	64	82	52,32	2.21e−2	30566	556	48 : 00 : 00
ste36a	9526	10966	8 992,22	9.72e−4	301900	130	48 : 00 : 00
ste36b	15852	18336	15 315,13	7.97e−4	301660	130	48 : 00 : 00
ste36c	8239110	9035686	7 980 802,09	8.86e−4	301558	130	48 : 00 : 00

Table 6

We solve SDP relaxations of QAP instances from TSPLIB using *SketchyCGAL*. We compute the relative gap and compare it with the values for the CSDP method [26] with clique size $k = \{2, 3, 4\}$ and the *PATH* method [116] reported in [26, Tab. 6]. Smaller is better.

Dataset	Optimum	CGAL	CSDP2	CSDP3	CSDP4	PATH
att48	10628	73.8	213	236.5	233.6	329.8
bayg29	1610	34.8	114.3	115.8	114.3	210.1
bays29	2020	55.7	107.6	118.3	115.4	164.8
berlin52	7542	52	175	127.2	127.2	280.6
bier127	118282	60.4	216.4	193.8	193.8	234.2
brazil58	25395	86.5	248	200.8	200.8	337
burma14	3323	10.8	24.6	28.4	32.3	95.5
ch130	6110	127.7	352.4	380.6	380.6	621.3
ch150	6528	121.1	346.9	318.2	318.2	689.3
dantzig42	699	73.5	193.1	174	174	82
eil101	629	70.6	227.3	235.3	235.3	437.7
eil51	426	57.3	203.6	205.4	205.5	244.4
eil76	538	65.4	282.9	183	183	328.2
fri26	937	32.7	91.6	39.4	39.4	41.6
gr120	6942	126.7	445.2	261.6	261.6	617.6
gr137	69853	147.5	264.6	220.3	220.3	38.9
gr17	2085	13.4	46.8	32.4	44.9	86.9
gr21	2707	27.2	94.5	69.7	66.3	185.7
gr24	1272	31.4	89.2	86.2	73.9	129.4
gr48	5046	59.7	210.2	187.4	187.4	270.4
gr96	55209	117.5	228.9	201.7	201.7	46
hk48	11461	43.9	222.4	207.7	207.7	281.6
kroA100	21282	125	469.6	469	469	720.2
kroA150	26524	184	411	467.4	467.4	945.8
kroB100	22141	156.8	411.9	313.6	313.6	624.2
kroB150	26130	151.9	417.3	353.7	353.7	844.7
kroC100	20749	169.5	507.4	445.1	445.1	763
kroD100	21294	118.5	504.2	349.8	349.8	654.4
kroE100	22068	152.3	489.5	346.3	346.3	684.2
lin105	14379	167.2	303.1	234.8	234.8	248.4
pr107	44303	190.5	181.5	207.9	207.9	41.6
pr124	59030	166.6	293.8	180.2	180.2	67.6
pr136	96772	144.5	325.5	164.7	164.7	196.6
pr144	58537	264.6	255	283.7	283.7	59.8
pr76	108159	81	192.2	194	194	39.4
rat99	1211	95.3	236.4	161.5	161.5	444.1
rd100	7910	88.8	438.4	375.3	375.3	506.5
st70	675	70.8	300.9	320	317.9	387.9
swiss42	1273	35.5	163.2	190.4	190.8	194
ulysses16	6859	11.1	23.6	20.2	23.2	82.7
ulysses22	7013	26.4	64.5	57	59.7	126.3

Table 7

We run *SketchyCGAL* for (the first of) 10^6 iterations or 48 hours of runtime, for solving SDP relaxations of QAP instances from TSPLIB. We report the upper bound after rounding (subsection 7.5.3), the objective value, the feasibility gap, the number of iterations, the memory usage (in MB), and the cpu time ('hh:mm:ss'). See subsection 7.5.

Dataset	Optimum	Upper bnd.	Objective	Feas. gap	Iteration	Memory	Time
att48	10628	18474	9 072,72	6.92e−4	311240	74	48 : 00 : 00
bayg29	1610	2170	1 498,60	5.40e−4	1000000	195	42 : 12 : 27
bays29	2020	3145	1 855,27	6.46e−4	1000000	195	42 : 19 : 21
berlin52	7542	11463	6 658,47	1.74e−3	266403	134	48 : 00 : 00
bier127	118282	189679	111 408,10	4.76e−2	23283	941	48 : 00 : 00
brazil58	25395	47362	18 336,86	1.70e−3	202400	129	48 : 00 : 00
burma14	3323	3682	3 153,63	5.14e−4	1000000	86	04 : 45 : 13
ch130	6110	13911	6 829,52	5.15e−2	20672	1077	48 : 00 : 00
ch150	6528	14432	11 021,38	1.39e−1	13925	1605	48 : 00 : 00
dantzig42	699	1213	589,35	5.11e−4	457221	129	48 : 00 : 00
eil101	629	1073	606,08	1.47e−2	46919	494	48 : 00 : 00
eil151	426	670	406,02	1.72e−3	286697	200	48 : 00 : 00
eil76	538	890	515,78	6.18e−3	103354	286	48 : 00 : 00
fri26	937	1243	858,28	6.57e−4	1000000	64	25 : 41 : 33
gr120	6942	15740	7 246,27	2.79e−2	27503	783	48 : 00 : 00
gr137	69853	172902	96 228,89	5.61e−2	19503	1206	48 : 00 : 00
gr17	2085	2365	1 800,49	3.83e−4	1000000	81	07 : 05 : 05
gr21	2707	3444	2 570,04	5.56e−4	1000000	64	12 : 12 : 17
gr24	1272	1672	1 136,65	3.53e−4	1000000	64	17 : 51 : 17
gr48	5046	8058	4 453,05	1.63e−3	357067	64	48 : 00 : 00
gr96	55209	120060	51 189,66	8.00e−3	58132	346	48 : 00 : 00
hk48	11461	16491	10 485,32	8.89e−4	357176	69	48 : 00 : 00
kroA100	21282	47891	19 953,20	8.72e−3	53552	462	48 : 00 : 00
kroA150	26524	75326	57 289,34	1.85e−1	11901	1605	48 : 00 : 00
kroB100	22141	56865	20 531,08	1.13e−2	43381	527	48 : 00 : 00
kroB150	26130	65832	57 292,86	2.02e−1	11722	1605	48 : 00 : 00
kroC100	20749	55920	19 831,30	1.08e−2	43352	509	48 : 00 : 00
kroD100	21294	46518	19 710,13	1.19e−2	43415	527	48 : 00 : 00
kroE100	22068	55680	20 463,81	1.06e−2	43352	557	48 : 00 : 00
lin105	14379	38417	12 640,70	1.13e−2	39662	440	48 : 00 : 00
pr107	44303	128712	36 333,09	1.32e−2	35009	651	48 : 00 : 00
pr124	59030	157381	74 152,15	3.76e−2	22049	900	48 : 00 : 00
pr136	96772	236593	135 053,09	6.06e−2	17446	1163	48 : 00 : 00
pr144	58537	213397	122 144,46	1.43e−1	14698	1383	48 : 00 : 00
pr76	108159	195718	91 294,27	2.50e−2	92634	287	48 : 00 : 00
rat99	1211	2365	1 206,81	9.76e−3	46602	462	48 : 00 : 00
rd100	7910	14935	7 459,54	9.89e−3	44863	527	48 : 00 : 00
st70	675	1153	586,71	2.92e−3	153971	134	48 : 00 : 00
swiss42	1273	1725	1 115,15	5.88e−4	599689	129	48 : 00 : 00
ulysses16	6859	7622	5 968,53	8.00e−4	1000000	73	05 : 16 : 03
ulysses22	7013	8865	5 999,82	1.03e−3	1000000	69	12 : 04 : 24

REFERENCES

- [1] P.-A. ABSIL, C. BAKER, AND K. GALLIVAN, *Trust-region methods on Riemannian manifolds*, Found. Comput. Math., 7 (2007), pp. 303–330.
- [2] A. Y. ALFAKIH, A. KHANDANI, AND H. WOLKOWICZ, *Solving Euclidean distance matrix completion problems via semidefinite programming*, Comput. Optim. Appl., 12 (1999), pp. 13–30.
- [3] F. ALIZADEH, *Combinatorial optimization with interior point methods and semi-definite matrices*, PhD thesis, Univ. Minnesota, 1991.
- [4] F. ALIZADEH, *Interior point methods in semidefinite programming with applications to combinatorial optimization*, SIAM J. Optim., 5 (1995), pp. 13–51.
- [5] F. ALIZADEH, J.-P. A. HAEBERLY, AND M. L. OVERTON, *Complementarity and nondegeneracy in semidefinite programming*, Math. Programming, 77 (1997), pp. 111–128.
- [6] F. ALIZADEH, J.-P. A. HAEBERLY, AND M. L. OVERTON, *Primal-dual interior point methods for semi-definite programming: Convergence rates, stability and numerical results*, SIAM J. Optim., 8 (1998), pp. 746–768.
- [7] Z. ALLEN-ZHU, Y. T. LEE, AND L. ORECCHIA, *Using optimization to obtain a width-independent, parallel, simpler, and faster positive SDP solver*, 2015, <https://arxiv.org/abs/1507.02259>.
- [8] Z. ALLEN-ZHU AND Y. LI, *Follow the compressed leader: Faster online learning of eigenvectors and faster MMWU*, June 2017, <https://arxiv.org/abs/1701.01722>.
- [9] M. F. ANJOS AND J. B. LASSERRE, *Handbook on Semidefinite, Conic and Polynomial Optimization*, Springer USA, 2012.
- [10] S. ARORA, E. HAZAN, AND S. KALE, *Fast algorithms for approximate semidefinite programming using the multiplicative weights update method*, in Proc. 46th Ann. IEEE Symp. Foundations of Computer Science, FOCS '05, Washington, DC, USA, 2005, pp. 339–348.
- [11] S. ARORA AND S. KALE, *A combinatorial, primal-dual approach to semidefinite programs*, J. ACM, 63 (2016), pp. 12:1–12:35.
- [12] F. BACH, *Duality between subgradient and conditional gradient methods*, SIAM J. Optim., 25 (2015), pp. 115–129.
- [13] D. BADER, H. MEYERHENKE, P. SANDERS, AND D. WAGNER, *Graph partitioning and graph clustering. 10th DIMACS Implementation Challenge Workshop*, February 2012, <https://www.cise.ufl.edu/research/sparse/matrices/DIMACS10/index.html> (accessed October 2019).
- [14] M. BAES, M. BÜRGISSE, AND A. NEMIROVSKI, *A randomized mirror-prox method for solving structured large-scale matrix saddle-point problems*, SIAM J. Optim., 23 (2013), pp. 934–962.
- [15] R. BALAN, B. G. BODMANN, P. G. CASAZZA, AND D. EDIDIN, *Painless reconstruction from magnitudes of frame coefficients*, J. Fourier Anal. Appl., 15 (2009), pp. 488–501.
- [16] R. BALAN, P. CASAZZA, AND D. EDIDIN, *On signal reconstruction without phase*, Appl. Comp. Harmonic Anal., 20 (2006), pp. 345–356.
- [17] A. S. BANDEIRA, *Convex relaxations for certain inverse problems on graphs*, PhD thesis, Princeton Univ., 2015.
- [18] A. BARVINOK, *A course in convexity*, AMS, Providence, RI, 2002.
- [19] A. BEN-TAL AND A. NEMIROVSKI, *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*, SIAM, Philadelphia, PA, 2001.
- [20] D. P. BERTSEKAS, *Constrained optimization and Lagrange multiplier methods*, Computer Science and Applied Mathematics, Academic Press, New York, NY, 1982.
- [21] S. BHOJANAPALLI, N. BOUMAL, P. JAIN, AND P. NETRAPALLI, *Smoothed analysis for low-rank solutions to semidefinite programs in quadratic penalty form*, in Proc. 31st Conf. Learning Theory, vol. 75, 2018, pp. 3243–3270.
- [22] N. BOUMAL, B. MISHRA, P.-A. ABSIL, AND R. SEPULCHRE, *Manopt, a Matlab toolbox for optimization on manifolds*, J. Mach. Learn. Res., 15 (2014), pp. 1455–1459.
- [23] N. BOUMAL, V. VORONINSKI, AND A. BANDEIRA, *The non-convex Burer–Monteiro approach works on smooth semidefinite programs*, in Adv. Neural Information Processing Systems 29, 2016, pp. 2757–2765.
- [24] S. BOYD, N. PARIKH, E. CHU, B. PELEATO, AND J. ECKSTEIN, *Distributed optimization and statistical learning via the alternating direction method of multipliers*, Found. Trends Mach. Learn., 3 (2011),

- pp. 1–122.
- [25] M. BRAND, *Fast low-rank modifications of the thin singular value decomposition*, Linear Alg. Appl., 415 (2006), pp. 20 – 30.
 - [26] J. F. S. BRAVO FERREIRA, Y. KHOO, AND A. SINGER, *Semidefinite programming approach for the quadratic assignment problem with a sparse graph*, Comput. Optim. Appl., 69 (2018), pp. 677–712.
 - [27] S. BURER AND R. D. MONTEIRO, *A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization*, Math. Prog., 95 (2003), pp. 329–357.
 - [28] S. BURER AND R. D. MONTEIRO, *Local minima and convergence in low-rank semidefinite programming*, Math. Prog., 103 (2005), pp. 427–444.
 - [29] S. BURER AND D. VANDENBUSSCHE, *Solving lift-and-project relaxations of binary integer programs*, SIAM J. Optim., 16 (2006), pp. 726–750.
 - [30] R. E. BURKARD, S. E. KARISCH, AND F. RENDL, *QAPLIB – a quadratic assignment problem library*, J. Global Opt., 10 (1997), pp. 391–403.
 - [31] E. J. CANDÈS, X. LI, AND M. SOLTANOLKOTABI, *Phase retrieval from coded diffraction patterns*, Appl. Comp. Harmonic Anal., 39 (2015), pp. 277 – 299.
 - [32] Y. CARMON, J. C. DUCHI, S. AARON, AND T. KEVIN, *A rank-1 sketch for matrix multiplicative weights*, in Proc. 32nd Conf. Learning Theory, vol. 99, 25–28 Jun 2019, pp. 589–623.
 - [33] A. CHAI, M. MOSCOSO, AND G. PAPANICOLAOU, *Array imaging using intensity-only measurements*, Inv. Prob., 27 (2010), p. 015005.
 - [34] D. CIFUENTES, *Burer–Monteiro guarantees for general semidefinite programs*, 2019, <https://arxiv.org/abs/1904.07147>.
 - [35] K. L. CLARKSON, *Coresets, sparse greedy approximation, and the Frank–Wolfe algorithm*, ACM Trans. Algorithms, 6 (2010), pp. 63:1–63:30.
 - [36] J. L. CONRAD, *Marburg virus hemorrhagic fever: Cytoarchitecture and histopathology*, <https://pixnio.com/de/wissenschaft/mikroskopische-aufnahmen/hamorrhagisches-fieber-marburg-virus/cytoarchitectural-histopathologischen-erkannt-niere-probe-marburg-patient> (accessed November 2019). Public domain (CC0 license).
 - [37] C. DELORME AND S. POLJAK, *Laplacian eigenvalues and the maximum cut problem*, Math. Prog., Ser. A, 62 (1993), pp. 557–574.
 - [38] C. DELORME AND S. POLJAK, *The performance of an eigenvalue bound on the max-cut problem in some classes of graphs*, Discrete Math., 111 (1993), pp. 145–156.
 - [39] L. DING, A. YURTSEVER, V. CEVHER, J. A. TROPP, AND M. UDELL, *An optimal-storage approach to semidefinite programming using approximate complementarity*, 2019, <https://arxiv.org/abs/1902.03373>.
 - [40] M. FAZEL, *Matrix Rank Minimization with Applications*, PhD thesis, Stanford University, Palo Alto, CA, USA, 2002.
 - [41] J. R. FIENUP, *Phase retrieval algorithms: a comparison*, Appl. Optics, 21 (1982), pp. 2758–2769.
 - [42] M. FRANK AND P. WOLFE, *An algorithm for quadratic programming*, Naval Res. Logistics Quart., 3 (1956), pp. 95–110.
 - [43] D. GARBER AND E. HAZAN, *Approximating semidefinite programs in sublinear time*, in Adv. Neural Information Processing Systems 25, 2011.
 - [44] D. GARBER AND E. HAZAN, *Sublinear time algorithms for approximate semidefinite programming*, Math. Prog., 158 (2016), pp. 329–361.
 - [45] G. GIDEL, F. PEDREGOSA, AND S. LACOSTE-JULIEN, *Frank–Wolfe splitting via augmented Lagrangian method*, in Proc. 21st Int. Conf. Artificial Intelligence and Statistics, vol. 84, 2018, pp. 1456–1465.
 - [46] A. GITTENS, *Topics in Randomized Numerical Linear Algebra*, PhD thesis, California Institute of Technology, Pasadena, CA, USA, 2013.
 - [47] M. X. GOEMANS AND D. P. WILLIAMSON, *Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming*, J. ACM, 42 (1995), pp. 1115–1145.
 - [48] G. H. GOLUB AND C. F. VAN LOAN, *Matrix computations*, Johns Hopkins Univ. Press, Baltimore, MD, fourth ed., 2013.
 - [49] N. HALKO, P. G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions*, SIAM Rev., 53 (2011), pp. 217–288.
 - [50] E. HAZAN, *Sparse approximate solutions to semidefinite programs*, in LATIN 2008: Theoretical Infor-

- maths, 2008, pp. 306–316.
- [51] C. HELMBERG, *Numerical evaluation of SBmethod*, Math. Prog., Ser. B, 95 (2003), pp. 381–406.
 - [52] C. HELMBERG AND F. OUSTRY, *Bundle methods to minimize the maximum eigenvalue function*, in Handbook of Semidefinite Programming, Springer, Boston, MA, 2000.
 - [53] C. HELMBERG AND F. RENDL, *A spectral bundle method for semidefinite programming*, SIAM J. Optim., 10 (2000), pp. 673–696.
 - [54] S. HOMER AND M. PEINADO, *Design and performance of parallel and distributed approximation algorithms for maxcut*, J. Parallel Distrib. Comput., 46 (1997), pp. 48 – 61.
 - [55] R. HORSTMAYER, R. Y. CHEN, X. OU, B. AMES, J. A. TROPP, AND C. YANG, *Solving ptychography with a convex relaxation*, New J. Phys., 17 (2015), p. 053044.
 - [56] Q. HUANG, Y. CHEN, AND L. GUIBAS, *Scalable semidefinite relaxation for maximum a posterior estimation*, in Proc. 31st Int. Conf. Machine Learning, vol. 32(2), 2014, pp. 64–72.
 - [57] M. JAGGI, *Revisiting Frank–Wolfe: Projection-free sparse convex optimization*, in Proc. 30th Int. Conf. Machine Learning, vol. 28(1), 2013, pp. 427–435.
 - [58] L. K. JONES, *A simple lemma on greedy approximation in Hilbert space and convergence rates for projection pursuit regression and neural network training*, Ann. Statist., 20 (1992), pp. 608–613.
 - [59] R. JONKER AND A. VOLGENANT, *A shortest augmenting path algorithm for dense and sparse linear assignment problems*, Computing, 38 (1987), pp. 325–340.
 - [60] A. P. KAMATH AND N. K. KARMARKAR, *A continuous method for computing bounds in integer quadratic optimization problems*, J. Global Optim., 2 (1991), pp. 229–241.
 - [61] A. P. KAMATH AND N. K. KARMARKAR, *An $O(nL)$ iteration algorithm for computing bounds in quadratic optimization problems*, in Complexity in Numerical Optimization, World Scientific Publishing, July 1993, pp. 254–268.
 - [62] N. K. KARMARKAR, *A new polynomial-time algorithm for linear programming*, Combinatorica, 4 (1984), pp. 373–395.
 - [63] R. M. KARP, *Reducibility among combinatorial problems*, in Complexity of computer computations, 1972, pp. 85–103.
 - [64] K. KRISHNAN AND T. TERLAKY, *Interior point and semidefinite approaches in combinatorial optimization*, in Graph Theory and Combinatorial Optimization, Springer, Boston, MA, 2005, ch. 1.
 - [65] J. KUCZYŃSKI AND H. WOŹNIAKOWSKI, *Estimating the largest eigenvalue by the power and Lanczos algorithms with a random start*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 1094–1122.
 - [66] H. KUHN, *The Hungarian Method for the assignment problem*, Naval Res. Logistics Quart., 2 (1955), pp. 83–97.
 - [67] B. KULIS, A. C. SURENDRAN, AND J. C. PLATT, *Fast low-rank semidefinite programming for embedding and clustering*, in Proc. 11th Int. Conf. Artificial Intelligence and Statistics, vol. 2, San Juan, Puerto Rico, 21–24 Mar 2007, pp. 235–242.
 - [68] J. LAVAEI AND S. H. LOW, *Zero duality gap in optimal power flow problem*, IEEE Trans. Power Sys., 27 (2012), pp. 92–107.
 - [69] Y. T. LEE AND S. PADMANABHAN, *An $\tilde{O}(m/\varepsilon^3.5)$ -cost algorithm for semidefinite programs with diagonal constraints*. Available at <http://arXiv.org/abs/1902.01859>, March 2019.
 - [70] R. B. LEHOUCQ, D. C. SORESENSEN, AND C. YANG, *ARPACK users’ guide*, vol. 6 of Software, Environments, and Tools, SIAM, Philadelphia, PA, 1998. Solution of large-scale eigenvalue problems with implicitly restarted Arnoldi methods.
 - [71] E. S. LEVITIN AND B. T. POLJAK, *Minimization methods in the presence of constraints*, Ž. Vychisl. Mat i Mat. Fiz., 6 (1966), pp. 787–823.
 - [72] H. LI, G. C. LINDERMAN, A. SZLAM, K. P. STANTON, Y. KLUGER, AND M. TYGERT, *Algorithm 971: an implementation of a randomized algorithm for principal component analysis*, ACM Trans. Math. Software, 43 (2017), pp. Art. 28, 14.
 - [73] Y.-F. LIU, X. LIU, AND S. MA, *On the nonergodic convergence rate of an inexact augmented lagrangian framework for composite convex programming*, Math. Oper. Res., 44 (2019), pp. 632–650.
 - [74] E. M. LOIOLA, N. M. M. DE ABREU, P. O. BOAVENTURA-NETTO, P. HAHN, AND T. QUERIDO, *A survey for the quadratic assignment problem*, Europ. J. Oper. Res., 176 (2007), pp. 657–690.
 - [75] A. MAJUMDAR, G. HALL, AND A. A. AHMADI, *A survey of recent scalability improvements for semidefinite programming with applications to machine learning, control, and robotics*. Available at

- <http://arXiv.org/abs/1908.05209>, Sep. 2019.
- [76] M. MESBAHI AND G. P. PAPAVALASSILOPOULOS, *On the rank minimization problem over a positive semi-definite linear matrix inequality*, IEEE Trans. Automatic Control, 42 (1997), pp. 239–243.
 - [77] L. MIRSKY, *Symmetric gauge functions and unitarily invariant norms*, Quart. J. Math. Oxford Ser. (2), 11 (1960), pp. 50–59.
 - [78] MOSEK APS, *The MOSEK optimization toolbox for MATLAB manual. Version 9.0.*, 2019.
 - [79] J. MUNKRES, *Algorithms for the assignment and transportation problems*, J. SIAM, 5 (1957), pp. 32–38.
 - [80] A. NEMIROVSKI, *Prox-method with rate of convergence $O(1/t)$ for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems*, SIAM J. Optim., 15 (2004), pp. 229–251.
 - [81] Y. NESTEROV, *Introductory Lectures on Convex Optimization: A Basic Course*, Springer US, 2004.
 - [82] Y. NESTEROV, *Primal-dual subgradient methods for convex problems*, Math. Program., 120 (2009), pp. 221–259.
 - [83] Y. NESTEROV AND A. NEMIROVSKI, *Self-concordant functions and polynomial time methods in convex programming*, Central Economic & Mathematic Institute report, USSR Acad. Sci., April 1990.
 - [84] Y. NESTEROV AND A. NEMIROVSKI, *Interior-Point Polynomial Algorithms in Convex Programming*, SIAM, Philadelphia, PA, 1994.
 - [85] M. L. OVERTON AND R. S. WOMERSLEY, *On the sum of the largest eigenvalues of a symmetric matrix*, SIAM J. Matrix Anal. Appl., 13 (1992), pp. 41–45, <https://doi.org/10.1137/0613006>.
 - [86] P. PARRILO, *Structured Semidefinite Programs and Semialgebraic Geometry Methods in Robustness and Optimization*, PhD thesis, California Institute of Technology, Pasadena, CA, 2000.
 - [87] G. PATAKI, *On the rank of extreme matrices in semidefinite programs and the multiplicity of optimal eigenvalues*, Math. Oper. Res., 23 (1998), pp. 339–358.
 - [88] R. PENG AND K. TANGWONGSAN, *Faster and simpler width-independent parallel algorithms for positive semidefinite programming*, in Proc. 24th Ann. ACM Symp. Parallelism in Algorithms and Architectures, 2012, pp. 101–108.
 - [89] T. PUMIR, S. JELASSI, AND N. BOUMAL, *Smoothed analysis of the low-rank approach for smooth semi-definite programs*, in Adv. Neural Information Processing Systems 31, 2018, pp. 2281–2290.
 - [90] A. RAGHUNATHAN, J. STEINHARDT, AND P. LIANG, *Semidefinite relaxations for certifying robustness to adversarial examples*, in Adv. Neural Information Processing Systems 32, 2018, pp. 10900–10910.
 - [91] G. REINELT, *TSPLIB*, <http://comopt.ifl.uni-heidelberg.de/software/TSPLIB95/>.
 - [92] D. M. ROSEN, L. CARLONE, A. S. BANDEIRA, AND J. J. LEONARD, *SE-sync: A certifiably correct algorithm for synchronization over the special euclidean group*, Int. J. Robot. Res., 38 (2019), pp. 95–125.
 - [93] M. F. SAHIN, A. EFTEKHARI, A. ALACAOGLU, F. LATORRE, AND V. CEVHER, *An inexact augmented Lagrangian framework for nonconvex optimization with nonlinear constraints*, 2019, <https://arxiv.org/abs/1906.11357>.
 - [94] S. SAHNI AND T. GONZALEZ, *P-complete approximation problems*, J. ACM, 23 (1976), pp. 555–565.
 - [95] A. SILVETI-FALLS, C. MOLINARI, AND J. FADILI, *Generalized conditional gradient with augmented Lagrangian for composite minimization*, 2019, <https://arxiv.org/abs/1901.01287>.
 - [96] N. SREBRO, *Learning with matrix factorizations*, PhD thesis, Massachusetts Institute of Technology, Cambridge, MA, 2004.
 - [97] J. STURM, *Using SeDuMi 1.02, a MATLAB toolbox for optimization over symmetric cones*, Optimization Methods and Software, 11–12 (1999), pp. 625–653.
 - [98] Y. SUN, Y. GUO, J. A. TROPP, AND M. UDELL, *Tensor random projection for low memory dimension reduction*, in NeurIPS Workshop on Relational Representation Learning, 2018.
 - [99] K. A. TOH, M. J. TODD, AND R. H. TÜTÜNCÜ, *SDPT3 — a Matlab software package for semidefinite programming, optimization methods and software*, Optim. Methods Software, 11 (1999), pp. 545–581.
 - [100] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Fixed-rank approximation of a positive-semidefinite matrix from streaming data*, in Adv. Neural Information Processing Systems 31, 2017, pp. 1225–1234.
 - [101] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Practical sketching algorithms for low-rank matrix approximation*, SIAM J. Matrix Anal. Appl., 38 (2017), pp. 1454–1485.
 - [102] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Streaming low-rank matrix approximation*

- with an application to scientific simulation, *SIAM J. Sci. Comput.*, 41 (2019), pp. A2430–A2463.
- [103] K. TSUDA, G. RÄTSCH, AND M. K. WARMUTH, *Matrix exponentiated gradient updates for on-line learning and Bregman projections*, *J. Mach. Learn. Res.*, 6 (2005), pp. 995–1018.
 - [104] I. WALDSPURGER, A. D’ASPREMONT, AND S. MALLAT, *Phase recovery, MaxCut and complex semidefinite programming*, *Math. Prog.*, 149 (2015), pp. 47–81.
 - [105] I. WALDSPURGER AND A. WATERS, *Rank optimality for the Burer–Monteiro factorization*, 2018, <https://arxiv.org/abs/1812.03046>.
 - [106] Z. WEN, D. GOLDFARB, S. MA, AND K. SCHEINBERG, *Row by row methods for semidefinite programming*, IEOR report, Columbia University, 2009.
 - [107] Z. WEN, D. GOLDFARB, AND W. YIN, *Alternating direction augmented Lagrangian methods for semidefinite programming*, *Math. Program. Comp.*, 2 (2010), pp. 203–230.
 - [108] L. YANG, D. SUN, AND K.-C. TOH, *SDPNAL+: A majorized semismooth Newton-CG augmented Lagrangian method for semidefinite programming with nonnegative constraints*, *Math. Prog. Comput.*, 7 (2015), pp. 331–366.
 - [109] Y. YE, *GSet random graphs*, <https://www.cise.ufl.edu/research/sparse/matrices/Gset/> (accessed October 2019).
 - [110] A. YURTSEVER, O. FERCOQ, AND V. CEVHER, *A conditional-gradient-based augmented Lagrangian framework*, in *Proc. 36th Int. Conf. Machine Learning*, vol. 97, 09–15 Jun 2019, pp. 7272–7281.
 - [111] A. YURTSEVER, O. FERCOQ, F. LOCATELLO, AND V. CEVHER, *A conditional gradient framework for composite convex minimization with applications to semidefinite programming*, in *Proc. 35th Intl. Conf. Machine Learning*, vol. 80, 2018, pp. 5727–5736.
 - [112] A. YURTSEVER, Y.-P. HSIEH, AND V. CEVHER, *Scalable convex methods for phase retrieval*, in *6th Int. IEEE Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, 2015, pp. 381–384.
 - [113] A. YURTSEVER, Q. TRAN DINH, AND V. CEVHER, *A universal primal-dual convex optimization framework*, in *Adv. Neural Information Processing Systems* 28, 2015, pp. 3150–3158.
 - [114] A. YURTSEVER, Q. TRAN-DINH, AND V. CEVHER, *Universal primal-dual proximal-gradient methods*, 2015, <https://arxiv.org/abs/1502.03123v2>.
 - [115] A. YURTSEVER, M. UDELL, J. TROPP, AND V. CEVHER, *Sketchy Decisions: Convex Low-Rank Matrix Optimization with Optimal Storage*, in *Proc. 20th Int. Conf. Artificial Intelligence and Statistics*, vol. 54, 20–22 Apr 2017, pp. 1188–1196.
 - [116] M. ZASLAVSKIY, F. BACH, AND J. VERT, *A path following algorithm for the graph matching problem*, *IEEE Trans. Pattern Anal. Mach. Intel.*, 31 (2009), pp. 2227–2242.
 - [117] Q. ZHAO, S. E. KARISCH, F. RENDL, AND H. WOLKOWICZ, *Semidefinite programming relaxations for the quadratic assignment problem*, *J. Comb. Optim.*, 2 (1998), pp. 71–109.
 - [118] X.-Y. ZHAO, D. SUN, AND K.-C. TOH, *A Newton-CG augmented Lagrangian method for semidefinite programming*, *SIAM J. Optim.*, 20 (2010), pp. 1737–1765.
 - [119] G. ZHENG, R. HORSTMAYER, AND C. YANG, *Wide-field, high-resolution Fourier ptychography microscopy*, *Nat. Photonics*, 7 (2013), pp. 739–745.